
Statistical Complexity of Soft Bellman Residual Minimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Bellman residual minimization (BRM) provides a scalable, gradient-based ap-
2 proach to offline reinforcement learning in large state spaces. While globally
3 convergent gradient-based soft (i.e., entropy-regularized) BRM methods for neural
4 networks have recently been established, their statistical complexity under
5 stochastic gradient descent remains largely unknown. In this paper, we address
6 this theoretical gap for soft BRM. Through a novel Lyapunov-based analysis, we
7 establish an $\mathcal{O}(1/n)$ average argument stability bound, which translates directly
8 into a $\mathcal{O}(1/n)$ statistical complexity for the soft BRM objective.

9 1 Introduction

10 Reinforcement learning (RL) provides a principled framework for sequential decision-making,
11 formalizing dynamic problems as Markov decision processes (MDPs) and learning policies that
12 optimize long-horizon rewards. In many high-stakes domains, however, online interaction is unethical,
13 dangerous, or prohibitively expensive [Jiang and Xie, 2024]. This motivates offline RL, where the
14 learner has access only to a fixed batch of data collected by a behavior policy, and its sister field,
15 inverse RL (IRL), which infers reward functions from logged expert trajectories.

16 Beyond the lack of online interaction, these problems pose both statistical and computational chal-
17 lenges: realistic state spaces are often high-dimensional or continuous, making exact dynamic
18 programming intractable. Following Jiang and Xie [2024], existing offline approaches can be broadly
19 organized into three paradigms: approximate dynamic programming, such as fitted Q -iteration
20 [Ernst et al., 2005]; marginalized importance sampling [Liu et al., 2018]; and Bellman Residual
21 Minimization (BRM) [Baird, 1995]. In this work, we adopt the BRM as our primary framework.

22 BRM reformulates dynamic programming as fitting a value function to the Bellman equation and
23 minimizes the mean-squared Bellman discrepancy under the offline data distribution. Despite its
24 conceptual appeal and empirical effectiveness, BRM has historically had limited theory for global
25 optimality and convergence. A notable exception is SBEED [Dai et al., 2018], which establishes
26 stable convergence guarantees for the soft, i.e., entropy-regularized, BRM objective, replacing the
27 Bellman max operator with a softmax¹.

28 Building on this direction, Kang et al. [2025] analyzed the optimization landscape of soft BRM.
29 After a classical bi-conjugate transformation, they showed that, for common Q -function classes
30 such as neural networks², the mean-squared Bellman error (MSBE) induces a Polyak–Łojasiewicz
31 (PL)–strongly-concave minimax landscape³. Consequently, plain stochastic gradient descent–ascent

¹A more comprehensive discussion of related work is deferred to Appendix D.

²Following Liu et al. [2022], this result uses neural networks whose width scales with $\text{radius}^{\text{depth}}$.

³The PL condition requires $\frac{1}{2}\|\nabla f(x)\|_2^2 \geq \mu(f(x) - f^*)$ for some $\mu > 0$, and yields convergence guarantees comparable to strong convexity even without convexity. In a PL–strongly-concave minimax problem, the inner maximization is strongly concave and the outer minimization satisfies the PL condition.

(SGDA) enjoys global convergence [Yang et al., 2020]. While this line of work substantially clarifies the optimization behavior of soft BRM, it leaves open the question of how statistical errors arising from finite offline data affect the learned policy:⁴ This raises the central question:

How many offline samples we need for BRM to recover a near-optimal value function, under SGDA?

Contributions. In this paper, we close this statistical gap for popular function classes such as neural networks and linear functions. Building on the PL structure identified by Kang et al. [2025], we develop a Lyapunov potential tailored to PL–strongly-concave optimization and blend it with a modern on-average argument-stability analysis. Our main theorem shows that, for a dataset of size n , SGDA attains an $\mathcal{O}(1/n)$ excess MSBE loss, a rate that doubles the exponent enjoyed in convex–concave optimization with Markov-sampled data, where the best known rate is $\mathcal{O}(1/\sqrt{n})$ [Wang et al., 2022]. In particular:

1. We prove $\mathcal{O}(1/n)$ on-average argument-stability bound for SGDA under PL–strong concavity, avoiding any independence assumptions on the minibatch sampling indices.
2. Leveraging this stability, we derive the $\mathcal{O}(1/n)$ generalization guarantee for BRM.
3. Our analysis is constructive, requires no variance reduction or extra regularisation, and applies verbatim to standard neural-network parameterisations commonly used in offline RL.

At a technical level, we couple two SGDA runs on neighboring datasets using the same initialization and the same minibatch index sequence. The Lyapunov potential contracts in expectation by a factor $(1 - c\eta_t)$, while the stochastic perturbations enter only through summable lower-order terms of order η_t^2 and η_t/n . Under Robbins–Monro stepsizes, the accumulated contraction dominates these perturbations, yielding the desired $\mathcal{O}(1/n)$ stability rate.

2 Setup and Backgrounds

Markov Decision Process. Throughout, we focus on a single-agent decision-making problem interacting with a discounted Markov Decision Process (MDP) described by the tuple $(\mathcal{S}, \mathcal{A}, P, r, \beta, \nu_0)$. A state is an element of the measurable space \mathcal{S} and the agent chooses actions from the finite set \mathcal{A} . For any state–action pair (s, a) the transition kernel $P(\cdot | s, a)$ gives a probability distribution over the next state and the reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ gives the immediate reward. We study the entropy-regularized control problem with regularization coefficient $\eta > 0$, where the Shannon entropy of a distribution $q \in \Delta_{\mathcal{A}}$ is $\mathcal{H}(q) := -\sum_{a \in \mathcal{A}} q(a) \log q(a)$. The scalar $\beta \in (0, 1)$ exponentially discounts rewards that occur further in the future and ν_0 denotes the distribution of the starting state s_0 . Throughout, the unsubscripted symbol η denotes the entropy-regularization coefficient; later, the subscripted sequence $(\eta_t)_{t \geq 0}$ continues to denote the SGDA stepsizes.

Policy and value functions. A (stationary Markov) policy $\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}$ assigns every state s to $\pi(\cdot | s)$, a distribution over actions \mathcal{A} ; when the agent is in state s_h at time h it samples $a_h \sim \pi(\cdot | s_h)$. Combined with the initial draw $s_0 \sim \nu_0$, a policy induces a probability measure $\mathbb{P}_{\nu_0, \pi}$ on infinite trajectories $(s_0, a_0, s_1, a_1, \dots)$, and the corresponding expectation operator is written $\mathbb{E}_{\nu_0, \pi}$. Under this setup, we consider the entropy-regularized optimal policy and its corresponding value functions defined as

$$\begin{aligned} \pi^* &:= \operatorname{argmax}_{\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}} \mathbb{E}_{\nu_0, \pi} \left[\sum_{h=0}^{\infty} \beta^h (r(s_h, a_h) + \eta \mathcal{H}(\pi(\cdot | s_h))) \right] \\ V^*(s) &:= \max_{\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}} \mathbb{E}_{\nu_0, \pi} \left[\sum_{h=0}^{\infty} \beta^h (r(s_h, a_h) + \eta \mathcal{H}(\pi(\cdot | s_h))) \mid s_0 = s \right] \\ Q^*(s, a) &:= \max_{\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}} \mathbb{E}_{\nu_0, \pi} \left[r(s_0, a_0) + \sum_{h=1}^{\infty} \beta^h (r(s_h, a_h) + \eta \mathcal{H}(\pi(\cdot | s_h))) \mid s_0 = s, a_0 = a \right]. \end{aligned}$$

⁴The main statistical guarantee of this paper was later referenced and used in the recent version of Kang et al. [2025]; see Appendix C.5 for details.

One can show that the optimal policy π^* and the value functions satisfy the following entropy-regularized Bellman optimality equations (Kang et al. [2025], Section 3.1 and Appendix B):

$$\begin{aligned} V^*(s) &= \max_{q \in \Delta_{\mathcal{A}}} \{ \mathbb{E}_{a \sim q} [Q^*(s, a)] + \eta \mathcal{H}(q) \} = \eta \ln \left[\sum_{a \in \mathcal{A}} \exp(Q^*(s, a)/\eta) \right] \\ \pi^*(a | s) &= \frac{\exp(Q^*(s, a)/\eta)}{\sum_{a' \in \mathcal{A}} \exp(Q^*(s, a')/\eta)} = \exp\left(\frac{Q^*(s, a) - V^*(s)}{\eta}\right) \text{ for } a \in \mathcal{A} \\ Q^*(s, a) &= r(s, a) + \beta \cdot \mathbb{E}_{s' \sim P(\cdot | s, a)} [V^*(s') | s, a] \end{aligned}$$

70 The normalized formulation recovered by setting $\eta = 1$ coincides with the mean-zero Gumbel-based
71 softmax Bellman equations used in the econometrics literature [Rust, 1994, Kang et al., 2025].

72 2.1 Bellman Residual Minimization

Bellman Error (Bellman Residual) and Temporal Difference Error. Define the function space \mathcal{Q} as the set of all bounded real-valued functions on the state-action space $\mathcal{S} \times \mathcal{A}$:

$$\mathcal{Q} := \{Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R} \mid \|Q\|_{\infty} < \infty\}$$

73 Under standard discounted-MDP assumption, the optimal action-value function, Q^* , is an element of
74 this space, i.e., $Q^* \in \mathcal{Q}$, provided that the discount factor β is in $(0, 1)$.

We introduce the (*soft*) *Bellman optimality operator*, $\mathcal{T} : \mathcal{Q} \rightarrow \mathcal{Q}$, defined for any $Q \in \mathcal{Q}$ by:

$$(\mathcal{T}Q)(s, a) := r(s, a) + \beta \cdot \mathbb{E}_{s' \sim P(\cdot | s, a)} [V_Q(s') | s, a] \text{ where } V_Q(s) := \eta \ln \left[\sum_{a \in \mathcal{A}} \exp(Q(s, a)/\eta) \right].$$

Note that Q^* is the unique fixed point of \mathcal{T} [Rust, 1994]; Q^* is the unique solution to the Bellman optimality equation:

$$\mathcal{T}Q^* = Q^* \text{ or equivalently, } (\mathcal{T}Q^*)(s, a) - Q^*(s, a) = 0.$$

75 The extent to which an arbitrary Q-function Q fails to satisfy the Bellman optimality equation
76 motivates the following definitions of error.

Definition 1 (Bellman Error (Bellman Residual)). *For any function $Q \in \mathcal{Q}$, we define the Bellman error (or Bellman residual) at a state-action pair (s, a) as the difference:*

$$(\mathcal{T}Q)(s, a) - Q(s, a)$$

Computing $\mathcal{T}Q$ requires knowledge of the transition kernel P , which is typically unavailable in reinforcement learning. Instead, given a sampled transition (s, a, s') , we work with the *sampled Bellman operator*, $\hat{\mathcal{T}}$. Given a single transition tuple (s, a, s') , this operator is defined as:

$$\hat{\mathcal{T}}Q(s, a, s') := r(s, a) + \beta \eta \log \sum_{a' \in \mathcal{A}} \exp(Q(s', a')/\eta)$$

Definition 2 (Temporal-Difference Error). *Using the sampled operator, we can define the Temporal-Difference (TD) error for a given transition (s, a, s') :*

$$\delta_Q(s, a, s') := \hat{\mathcal{T}}Q(s, a, s') - Q(s, a)$$

77 The connection between the Bellman error and the TD error is established in the following lemma. It
78 shows that the TD error is an unbiased, single-sample estimate of the Bellman error.

Lemma 1 (Relationship between Bellman and TD Errors). *For any $Q \in \mathcal{Q}$ and any state-action pair (s, a) , the expectation of the Sampled Bellman operator over the next state s' recovers the original Bellman operator:*

$$\mathbb{E}_{s' \sim P(\cdot | s, a)} [\hat{\mathcal{T}}Q(s, a, s')] = (\mathcal{T}Q)(s, a)$$

Consequently, the expected TD error is equal to the Bellman error:

$$\mathbb{E}_{s' \sim P(\cdot | s, a)} [\delta_Q(s, a, s')] = (\mathcal{T}Q)(s, a) - Q(s, a)$$

Bellman Residual Minimization. Note that both Bellman error (Bellman residual) and its proxy, the TD error, are functions of (s, a) . To find Q that minimizes the Bellman error for all (s, a) , we can instead find Q that minimizes expected square error on the offline data distribution. That is, we first define the *Squared Bellman Error* at (s, a) as $\mathcal{L}_{\text{BE}}(Q)(s, a) := ((\mathcal{T}Q)(s, a) - Q(s, a))^2$ and minimize the *Mean Squared Bellman Error* (MSBE), defined as:

$$\overline{\mathcal{L}}_{\text{BE}}(Q) := \mathbb{E}_{(s,a) \sim \pi_D, \nu_0} [\mathcal{L}_{\text{BE}}(Q)(s, a)]$$

where π_D is the policy used for collecting data. Furthermore, as a proxy for Squared Bellman Error, we define the *Squared TD Error*: $\mathcal{L}_{\text{TD}}(Q)(s, a, s') := \delta_Q(s, a, s')^2$ and minimize the *Mean Squared TD Error* (MSTDE) as a proxy for MSBE, defined as:

$$\overline{\mathcal{L}}_{\text{TD}}(Q) := \mathbb{E}_{(s,a) \sim \pi_D, \nu_0} [\mathbb{E}_{s' \sim P(\cdot|s,a)} [\mathcal{L}_{\text{TD}}(Q)(s, a, s')]]$$

79 Unfortunately, MSTDE is a *biased* proxy for MSBE. This bias happens because
80 $\mathbb{E}_{s' \sim P(\cdot|s,a)} [\delta_Q(s, a, s')^2 | s, a]^2 \neq \mathbb{E}_{s' \sim P(\cdot|s,a)} [\delta_Q(s, a, s')^2 | s, a]$ i.e., expectation and square
81 are not exchangeable. This issue is often called the *double sampling problem* [Antos et al., 2008].
82 Specifically, one can show that

$$\begin{aligned} \mathcal{L}_{\text{BE}}(Q)(s, a) &= \mathbb{E}_{s' \sim P(\cdot|s,a)} [\mathcal{L}_{\text{TD}}(Q)(s, a, s')] - \mathbb{E}_{s' \sim P(\cdot|s,a)} \left[\left(\mathcal{T}Q(s, a) - \hat{\mathcal{T}}Q(s, a, s') \right)^2 \right] \\ &= \mathbb{E}_{s' \sim P(\cdot|s,a)} [\mathcal{L}_{\text{TD}}(Q)(s, a, s')] - \beta^2 \mathbb{E}_{s' \sim P(\cdot|s,a)} \left[(V_Q(s') - \mathbb{E}_{s' \sim P(\cdot|s,a)} [V_Q(s')])^2 \right]. \end{aligned}$$

(For the detailed derivation, see [Kang et al., 2025, Appendix C.1].) Since the bias term includes the $\mathbb{E}_{s' \sim P(\cdot|s,a)}$ part, correcting this bias term again remains challenging without full knowledge of the system’s transition dynamics, P . To resolve this issue, we employ an approach often referred to as the “Bi-Conjugate Trick” [Antos et al., 2008, Dai et al., 2018, Patterson et al., 2022]:

$$\mathcal{L}_{\text{BE}}(Q)(s, a) = \mathbb{E}_{s' \sim P(\cdot|s,a)} [\delta_Q(s, a, s') | s, a]^2 = \max_{h \in \mathbb{R}} 2 \cdot \mathbb{E}_{s' \sim P(\cdot|s,a)} [\delta_Q(s, a, s') | s, a] \cdot h - h^2$$

83 According to [Kang et al., 2025, Appendix C.1], this bi-conjugate form can be re-parametrized using
84 $\zeta := h - r(s, a) + Q(s, a)$ as:

$$\mathcal{L}_{\text{BE}}(Q)(s, a) = \mathbb{E}_{s' \sim P(\cdot|s,a)} [\mathcal{L}_{\text{TD}}(Q)(s, a, s')] - \beta^2 \min_{\zeta \in \mathbb{R}} \mathbb{E}_{s' \sim P(\cdot|s,a)} [(V_Q(s') - \zeta)^2] \quad (1)$$

85 Now define the per-sample objective

$$f(Q, \zeta; s, a, s') := \mathcal{L}_{\text{TD}}(Q)(s, a, s') - \beta^2 (V_Q(s') - \zeta)^2. \quad (2)$$

86 If we slightly abuse notation such that $\mathbb{E}_{s,a,s'}[*] = \mathbb{E}_{s,a \sim \nu_0, \pi_D} [\mathbb{E}_{s' \sim P(\cdot|s,a)}[*]]$, then minimizing
87 MSBE can be written as the following min-max problem:

$$\min_{Q \in \mathcal{Q}} \mathbb{E}_{(s,a,s')} [f(Q, \tilde{\zeta}(Q); s, a, s')] \quad \text{where } \tilde{\zeta}(Q) \in \operatorname{argmin}_{\zeta \in \mathbb{R}^{S \times A}} \mathbb{E}_{(s,a,s')} [(V_Q(s') - \zeta(s, a))^2].$$

88 By parametrizing ζ as ζ_v and Q as Q_w using function classes such as neural networks and relaxing
89 the definition of the per-sample objective

$$f(w, v; s, a, s') := f(Q_w, \zeta_v; s, a, s') = \mathcal{L}_{\text{TD}}(Q_w)(s, a, s') - \beta^2 (V_{Q_w}(s') - \zeta_v)^2, \quad (3)$$

90 the resulting parametrized problem becomes

$$\min_{w \in \mathcal{W}} \mathbb{E}_{(s,a,s')} [f(w, \tilde{v}; s, a, s')] \quad \text{where } \tilde{v}(w) \in \operatorname{argmin}_{v \in \mathcal{V}} \mathbb{E}_{(s,a,s')} [(V_{Q_w}(s') - \zeta_v(s, a))^2]. \quad (4)$$

91 Let’s consider Equation (4) as the expected risk minimization problem. Define $\mathcal{D}_{\pi_D, \nu_0} = \{z_i\}_{i=1}^N$ as
92 the offline data collected from following π_D starting from ν_0 , where each $z_i = (s_i, a_i, s'_i)$ denotes
93 a sample consisting of the current state, the action taken, and the resulting next state. Then, the
94 corresponding empirical risk minimization problem can be written as

$$\min_{w \in \mathcal{W}} \frac{1}{N} \sum_{z \in \mathcal{D}_{\pi_D, \nu_0}} f(w, \tilde{v}; z), \quad \text{where } \tilde{v}(w) \in \operatorname{argmin}_{v \in \mathcal{V}} \frac{1}{N} \sum_{z \in \mathcal{D}_{\pi_D, \nu_0}} [(V_{Q_w}(s') - \zeta_v(s, a))^2]. \quad (5)$$

95 A canonical way of solving the mini-max problem is to apply the Stochastic Gradient Ascent Descent
96 (SGDA) algorithm [Yang et al., 2020]. Kang et al. [2025] proved that both (4) and (5) satisfy the
97 Polyak-Łojasiewicz condition within a large enough ball around the initialization point for a Neural
98 Network parametrization of Q with a sufficient network width (width scaling with radius^{depth}), and
99 therefore the SGDA algorithm finds the (w, v) that are global minima of Equation (5). (For the
100 detailed discussion on neural networks, see Appendix C.2.)

101 **2.2 Stochastic Gradient Ascent–Descent Algorithm (SGDA)**

102 As discussed earlier, Stochastic Gradient Ascent–Descent (SGDA) is the workhorse we use to solve
 103 the minimax problem (5). Given a function $f(w, v)$, at every iteration it performs a *descent* step
 104 on the primal variable w and an *ascent* step on the dual variable v using (possibly noisy) gradients
 105 computed from a minibatch of samples.

106 Let $\mathcal{D} = \{z_i\}_{i=1}^n$ denote the dataset, where each $z_i = (s_i, a_i, s'_i)$ denotes a sample consisting of the
 107 current state, the action taken, and the resulting next state. Fix a minibatch size $B \in \{1, \dots, n\}$. At
 108 round t we draw an index set $I_t \subseteq [n]$, $|I_t| = B$, either *with* or *without* replacement (our theory does
 109 not depend on this choice). With f standing for the relaxed per–sample saddle objective introduced
 110 in (3), the averaged stochastic gradients are

$$g_t^w = \nabla_w F_{I_t}(w_t, v_t), \quad g_t^v = \nabla_v F_{I_t}(w_t, v_t) \quad \text{where } F_I(w, v) := \frac{1}{|I|} \sum_{j \in I_t} f(w, v; z_j) \text{ for any } I \subset [n].$$

111 By a slight abuse of notation, when the subscript corresponds to the full dataset \mathcal{D} , we write
 112 $F_{\mathcal{D}}(w, v) := \frac{1}{n} \sum_{z \in \mathcal{D}} f(w, v; z)$.

113 Unbiasedness is preserved: $\mathbb{E}[g_t^w] = \nabla_w F_{\mathcal{D}}(w_t, v_t)$ and likewise for v , while the variance contracts
 114 by the usual $1/B$ factor. Using stepsize sequence $(\eta_t)_{t \geq 0}$, SGDA proceeds as

$$w_{t+1} = w_t - \eta_t g_t^w, \quad v_{t+1} = v_t + \eta_t g_t^v.$$

115 The recursion can be written as follows, which is the form used in the stability proofs of Section 3.

$$(w_{t+1}, v_{t+1}) = (w_t, v_t) + \eta_t (-g_t^w, g_t^v),$$

Algorithm 1 Minibatch SGDA on the empirical objective (5)

Input: Dataset $\mathcal{D} = \{z_i\}_{i=1}^n$, minibatch size B , stepsizes (η_t) , initial (w_0, v_0)

for $t = 0$ **to** $T - 1$ **do**

 Draw $I_t \subseteq [n]$ with $|I_t| = B$ uniformly at random

$g_t^w \leftarrow \frac{1}{B} \sum_{i \in I_t} \nabla_w f(w_t, v_t; z_i)$

$g_t^v \leftarrow \frac{1}{B} \sum_{i \in I_t} \nabla_v f(w_t, v_t; z_i)$

$w_{t+1} \leftarrow w_t - \eta_t g_t^w$ {gradient descent}

$v_{t+1} \leftarrow v_t + \eta_t g_t^v$ {gradient ascent}

end for

Output: (w_T, v_T)

116 With harmonic–stepsizes, we can state the global convergence guarantee of ALGORITHM 1 in the
 117 parameter space:

118 **Lemma 2** (Global convergence of minibatch SGDA in parameter space [Yang et al., 2020, Kang
 119 et al., 2025]). *Let the iterates in ALGORITHM 1 be written as $\{(w_t, v_t)\}_{t \geq 0}$, where w_t parametrises
 120 the action–value function Q_{w_t} (primal variable) and v_t parametrises ζ_{v_t} (dual variable). Choose
 121 the harmonic step-sizes $\eta_t = \frac{c_1}{c_2 + t}$, $t \geq 0$ with some constants $c_1 > 0$ and $c_2 \geq 1$ such that
 122 $\eta_t \leq \min\{1/(4L), 1/\rho\}$ for every t . Then the ALGORITHM 1’s output sequence $\{(w_t, v_t)\}$ converges
 123 almost surely to a saddle point of the empirical objective (5), where the suboptimality of empirical
 124 objective (5) is bounded by $\frac{d_1}{d_2 + t}$ for some constants d_1 and d_2 .*

125 Lemma 2 clarifies the optimization side of empirical BRM: under the PL–strongly-concave land-
 126 scape, SGDA globally converges to a saddle point of the empirical objective. What remains is the
 127 generalization side. The lemma does not address how errors induced by the finite offline dataset
 128 propagate to the population BRM objective, or whether the value function learned from empirical
 129 data generalizes to the underlying MDP distribution.

130 **3 Stability and Generalization for Bellman Residual Minimization**

131 We quantify generalization through *algorithmic stability* for minimax learning. Algorithmic stability
 132 formalizes how sensitive a learning algorithm is to small changes in the training set: if replacing a

133 single training example only slightly perturbs the algorithm’s output, then the algorithm is said to be
 134 stable. The key fact is that stability implies generalization: algorithms that are stable on neighbouring
 135 datasets exhibit small discrepancies between their empirical risk and population risk. As shown in
 136 Wang et al. [2022], this principle carries over to minimax optimization.

137 **Notations.** For the stability analysis, it is convenient to momentarily step back from the specific
 138 Bellman–residual objective and view our problem as a generic empirical PL–strong-concave saddle-
 139 point problem. In the BRM formulation, the primal variable w parametrizes the action–value function
 140 Q_w , the dual variable v parametrizes the auxiliary function ζ_v , and the empirical objective in (5) is
 141 an average over samples $z = (s, a, s')$ drawn from the offline dataset. In this section we abstract
 142 this structure and write $f(w, v; z)$ for the per–sample saddle loss and $F_D(w, v)$ for its empirical
 143 average over a dataset $D = \{z_i\}_{i=1}^n$. Working in this slightly more abstract template keeps the proofs
 144 uncluttered; in Theorem 6 we then specialize the resulting stability and generalization bounds back to
 145 BRM.

146 In this section, for the sake of generality, we use notations that generalize the Bellman residual
 147 minimization problem. Let $\mathcal{D} = (z_1, \dots, z_n)$ be a dataset with $z_i \in \mathcal{Z}$. For any $i \in [n]$ and any
 148 $\tilde{z}_i \in \mathcal{Z}$, define the *replace-one neighbour* of \mathcal{D} by $\mathcal{D}^{(i)} := (z_1, \dots, z_{i-1}, \tilde{z}_i, z_{i+1}, \dots, z_n)$. We call
 149 \mathcal{D} and $\mathcal{D}^{(i)}$ *neighbouring datasets*. When comparing SGDA runs on \mathcal{D} and $\mathcal{D}^{(i)}$, we couple them
 150 using the same minibatch index sequence $(i_t)_{t \geq 0}$ (shared-index coupling). Expectations averaged
 151 over i are taken with $i \sim \text{Unif}([n])$. This coupling is the key mechanism that allows the analysis to
 152 proceed without an i.i.d. assumption on the data. To see how it works, please check Appendix C.3.

153 Throughout, we present all proofs for the *single-sample* (“minibatch-of-one”) variant of SGDA. The
 154 extension to a minibatch of size $B \geq 1$ (or the full-batch case. $B = n$) is mechanical. Every stochastic-
 155 gradient term is replaced by its averaged counterpart, which reduces all variance contributions by a
 156 factor $1/B$, while the probability that a particular data point appears in the update increases from
 157 $1/n$ to B/n . Consequently, every lemma and theorem below remains valid verbatim—with constants
 158 rescaled by these factors—and no new conceptual issues arise.

159 3.1 Stability

160 We consider an offline setting where the data may be dependently sampled (e.g., from a single
 161 trajectory in a Markov Decision Process), violating the standard i.i.d. assumption. In this case, as in
 162 Wang et al. [2022], the concept of *on-average algorithmic stability* is useful.

163 **Definition 3** (On-average algorithmic stability). *Let \mathcal{A} be a randomized learning algorithm that*
 164 *maps a dataset $\mathcal{D} = (z_1, \dots, z_n)$ to a parameter output $\mathcal{A}(\mathcal{D})$. For each $i \in [n]$, let $\mathcal{D}^{(i)}$ denote the*
 165 *replace-one neighbor of \mathcal{D} , where z_i is replaced by an independent copy \tilde{z}_i . Then the on-average*
 166 *argument stability of \mathcal{A} after T iterations is $\varepsilon_T := \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\|\mathcal{A}(\mathcal{D}) - \mathcal{A}(\mathcal{D}^{(i)})\|]$, where the*
 167 *expectation is over the algorithm’s internal randomness and the choice of i .*

168 Intuitively, ε_T measures how much the algorithm’s output changes when a single training point is
 169 replaced on average. We analyze the on-average argument stability of SGDA (i.e., $\mathcal{A} = \text{SGDA}$)
 170 for smooth–strongly concave saddle problems, following the stability framework for minimax
 171 optimization.

172 Let SGDA iterates start from the same initialization $(w_0, v_0) = (w'_0, v'_0)$:

$$\begin{aligned} w_{t+1} &:= w_t - \eta_t \nabla_w f(w_t, v_t; z_{i_t}), & w'_{t+1} &:= w'_t - \eta_t \nabla_w f(w'_t, v'_t; z'_{i_t}), \\ v_{t+1} &:= v_t + \eta_t \nabla_v f(w_t, v_t; z_{i_t}), & v'_{t+1} &:= v'_t + \eta_t \nabla_v f(w'_t, v'_t; z'_{i_t}), \end{aligned}$$

173 with same-index coupling of datasets: $z'_j = z_j$ for $j \neq i$ and $z'_i = \tilde{z}_i$. For $D \in \{\mathcal{D}, \mathcal{D}^{(i)}\}$, define

$$F_D(w, v) := \frac{1}{n} \sum_{j=1}^n f(w, v; z_j^D), \quad \Phi_D(w) := \max_v F_D(w, v), \quad \Phi_D^* := \min_w \Phi_D(w).$$

174 We let $\mathcal{F}_t := \sigma((w_s, v_s, w'_s, v'_s, i_s)_{0 \leq s \leq t})$ be the natural filtration and introduce a *ghost* index
 175 $\hat{i}_t \sim \text{Unif}(\{1, \dots, n\})$ independent of \mathcal{F}_t , shared by both runs. The role of the ghost index is to
 176 decouple the sampling noise at time t from the past and from the coupling across datasets.

177 Throughout, we impose the following conditions (C1)–(C6) satisfied by the BRM objective [Kang
178 et al., 2025] under Section 2 when the neural network parametrization’s width is enough⁵:

- 179 (C1) **Smoothness.** F_D is L -smooth in the joint variable (w, v) .
180 (C2) **PL for Φ_D .** Φ_D satisfies the Polyak–Łojasiewicz (PL) inequality with parameter $\mu_{\text{PL}} > 0$:
181 $\frac{1}{2} \|\nabla \Phi_D(w)\|^2 \geq \mu_{\text{PL}} (\Phi_D(w) - \Phi_D^*)$.
182 (C3) **QG for Φ_D .** Φ_D satisfies Quadratic Growth (QG) with parameter $\mu_{\text{QG}} > 0$: $\Phi_D(w) - \Phi_D^* \geq$
183 $\frac{\mu_{\text{QG}}}{2} \|w - x_D^*\|^2$.
184 (C4) **Strong concavity in v .** $F_D(\cdot, \cdot)$ is ρ -strongly concave in v uniformly in w .
185 (C5) **Bounded gradients on the effective domain.** There exists a compact convex set $\Omega \subset \mathcal{W} \times \mathcal{V}$
186 such that the sequence of iterates $\{(w_t, v_t)\}_{t=0}^T$ generated by the algorithm remains within Ω
187 almost surely. We define $G < \infty$ as the uniform gradient bound on this set:

$$G := \sup_{(w,v) \in \Omega, z \in \mathcal{Z}} \|\nabla_w f(w, v; z)\| \vee \|\nabla_v f(w, v; z)\|.$$

- 188 (C6) **Selected saddle point.** Each dataset D admits at least one saddle point, and the SGDA iterates
189 initialized at (w_0, v_0) converge to an initialization-dependent saddle point (x_D^*, v_D^*) , which we
190 call the selected saddle.

191 In addition, we impose the following settings for the stability analysis. (S1) is a relatively trivial
192 setting, (S2) is the key mechanism that allows the analysis to proceed without an i.i.d. assumption on
193 the data, and (S3) is standard in the stability literature: for more details, see Appendix C.4.

- 194 (S1) **Stepsizes.** $0 < \eta_t \leq \min\{\frac{1}{4L}, \frac{1}{\rho}\}$.
195 (S2) **Shared-index coupling (no i.i.d. needed).** The two coupled runs on \mathcal{D} and $\mathcal{D}^{(i)}$ use the *same*
196 index sequence $(i_t)_{t \geq 0}$.
197 (S3) **Uniformity of constants across datasets.** The constants $L, \rho, \mu_{\text{PL}}, \mu_{\text{QG}}, G$ are the same for
198 D and $D^{(i)}$.

199 To establish stability in our PL–strongly-concave setting, a key difficulty is that SGDA is not
200 optimizing the value function $\Phi_D(w) := \max_v F_D(w, v)$ directly, but rather the saddle objective
201 $F_D(w, v)$. The gradient used for the primal update, $\nabla_w F_D(w_t, v_t)$, coincides with $\nabla \Phi_D(w_t)$ only
202 when the dual variable v_t is already at its inner maximizer $v_D^*(w_t)$; away from this manifold there
203 is a non-negligible “mismatch” term $\Delta_t := \nabla_w F_D(w_t, v_t) - \nabla \Phi_D(w_t)$. As a result, tracking the
204 PL suboptimality $\Phi_D(w_t) - \Phi_D^*$ alone does not yield a contracting recursion under SGDA, while
205 tracking only the dual gap $\Phi_D(w_t) - F_D(w_t, v_t)$ ignores how far the primal iterate is from the BRM
206 minimizer.

207 To address this, we define the *Lyapunov potential* :

$$\Psi_{\alpha, D}(w, v) := \underbrace{\Phi_D(w) - \Phi_D^*}_{A(w)} + \alpha \cdot \underbrace{(\Phi_D(w) - F_D(w, v))}_{\Gamma(w, v)},$$

208 where $\alpha \in [\frac{4L^2}{\rho^2}, \infty)$. Our Lyapunov potential is designed precisely to couple these two effects
209 in a single scalar quantity: the PL term measures how close the current value function is to the
210 data-dependent BRM optimum, and the dual-gap term penalizes the inner mismatch strongly enough
211 that the beneficial drift from strong concavity in v dominates the adverse contribution of Δ_t to the
212 primal dynamics.

213 With a suitable choice of the weight α , this potential contracts in expectation at each SGDA step, and
214 can then be translated back into a bound on the distance between coupled SGDA trajectories via PL
215 and QG. Because the step sizes in Lemma 2 satisfy the Robbins–Monro conditions $\sum_t \eta_t = \infty$ and
216 $\sum_t \eta_t^2 < \infty$, the contraction dominates these perturbations and keeps the two trajectories close for
217 large t . The next theorem formalizes this as an $\mathcal{O}(1/n)$ bound on the on-average argument stability
218 of SGDA, stated in a general PL–strongly-concave minimax setting that our BRM formulation
219 instantiates.

⁵For details, See Appendix C.2.

220 **Theorem 3** (On-average argument stability of SGDA without i.i.d. sampling). *Suppose that the*
 221 *loss function f and datasets $\mathcal{D}, \mathcal{D}^{(i)}$ satisfy (C1)–(C6) and under the setting (S1)–(S3). Let $\varepsilon_T :=$*
 222 *$\frac{1}{n} \sum_{i=1}^n \mathbb{E}[\|w_T(\mathcal{D}) - w_T(\mathcal{D}^{(i)})\| + \|v_T(\mathcal{D}) - v_T(\mathcal{D}^{(i)})\|]$. With choices of $\beta := 2L$ and α above,*
 223 *the following on-average argument stability bound holds:⁶*

$$\varepsilon_T \leq 2 C_{\text{dist}} \left[e^{-\frac{3c}{4} \sum_{s=0}^{T-1} \eta_s} \Psi_{\alpha,0}^{\max} + C'_{\text{var}} \sum_{t=0}^{T-1} \eta_t^2 e^{-\frac{3c}{4} \sum_{s=t+1}^{T-1} \eta_s} \right]^{1/2} + \frac{C'_{\text{hit}}}{n}$$

224 where $\Psi_{\alpha,0}^{\max} := \max\{\mathbb{E}[\Psi_{\alpha,\mathcal{D}}(w_0, v_0)], \max_{1 \leq i \leq n} \mathbb{E}[\Psi_{\alpha,\mathcal{D}^{(i)}}(w_0, v_0)]\}$, and $c, C'_{\text{var}}, C_{\text{dist}}, C'_{\text{hit}}$
 225 are all constants depending only on $L, \rho, \mu_{\text{PL}}, \mu_{\text{QG}}, G$.⁷

226 To make the rate transparent, we now specialize the stepsizes to the harmonic Robbins–Monro rule
 227 $\eta_t = c_1/(c_2 + t)$, which Kang et al. [2025] chooses for Algorithm 1 to prove Lemma 2. This schedule
 228 satisfies Setting (S1) for suitable $c_1 > 0$, $c_2 \geq 1$ and turns the kernel sums in Theorem 3 into closed
 229 forms. This yields the next corollary, which displays roughly $O(T^{-1/2})$ decay of the optimization
 230 term while keeping the $O(1/n)$ contribution explicit. See Appendix B.2 for details.

231 **Corollary 4** (Informal bound under a harmonic stepsize schedule). *Choose the Robbins–Monro*
 232 *stepsizes in ALGORITHM 1 as $\eta_t = \frac{c_1}{c_2 + t}$, $t \geq 0$, with constants $c_1 > 0$ and $c_2 \geq 1$ small enough*
 233 *that satisfying (S1). Let $c := \min\{\mu_{\text{PL}}/2, \rho/2\}$. Then from Theorem 3, the stability constant ε_T in*
 234 *Theorem 3 admits the stability bound scales as*

$$\varepsilon_T = O\left((c_2 + T)^{-\min\{\frac{1}{2}, \frac{3c c_1}{8}\}}\right) + O\left(\frac{1}{n}\right).$$

235 3.2 Generalization

236 In this section, we quantify generalization through algorithmic stability we derived in the previous
 237 section. Following the minimax stability framework of Wang et al. [2022], stability controls both
 238 the *primal function*, which is the Bellman residual, and the *weak primal–dual gap*. We first define
 239 these two risks, then invoke the transfer lemma that turns our stability bound from Theorem 3 into
 240 generalization guarantees. Specifically, our goal is to 1) bound the difference between population Bell-
 241 man–residual risk and empirical Bellman–residual risk and 2) bound the population Bellman–residual
 242 risk of the SGDA output.

243 **Definition 4** (Primal Risk). *Under (C4), define the value function*

$$R(w) := \max_{v \in \mathcal{V}} F(w, v), \quad R_n(w) := \max_{v \in \mathcal{V}} F_{\mathcal{D}}(w, v).$$

244 **Definition 5** (Weak primal–dual risk). *For $(w, v) \in \mathcal{W} \times \mathcal{V}$, define the population and empirical*
 245 *weak–PD risks by*

$$\Delta^{\text{PD}}(w, v) := \max_{v' \in \mathcal{V}} F(w, v') - \min_{w' \in \mathcal{W}} F(w', v), \quad \Delta_n^{\text{PD}}(w, v) := \max_{v' \in \mathcal{V}} F_{\mathcal{D}}(w, v') - \min_{w' \in \mathcal{W}} F_{\mathcal{D}}(w', v).$$

246 The key transfer principle is stability \Rightarrow generalization for minimax problems: if the SGDA iterate
 247 (w_T, v_T) has on-average argument stability ε_T on neighboring datasets, then the primal value–function
 248 gap $\mathbb{E}[R(w_T) - R_n(w_T)]$ and the weak primal–dual gap $|\mathbb{E}[\Delta^{\text{PD}}(w_T, v_T) - \Delta_n^{\text{PD}}(w_T, v_T)]|$ can be
 249 effectively upper bounded by constant times ε_T from Theorem 3, which is $\mathcal{O}(1/n)$. The remainder
 250 of this subsection formalizes this to apply this transfer with our stability bound (Theorem 3) to obtain
 251 Theorem 6, $\mathcal{O}(1/n)$ generalization for BRM under SGDA.

252 **Lemma 5** (Theorem 5, Wang et al. [2022]). *Let Conditions (C1), (C4), and (C5) hold. Let (w_T, v_T)*
 253 *be the SGDA iterates produced on \mathcal{D}_n and let*

$$\varepsilon_T = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\|w_T(\mathcal{D}) - w_T(\mathcal{D}^{(i)})\| + \|v_T(\mathcal{D}) - v_T(\mathcal{D}^{(i)})\|]$$

⁶Note that under the Robbins–Monro conditions $\sum_t \eta_t = \infty$ and $\sum_t \eta_t^2 < \infty$, the optimization term (the square root term) vanishes as $T \rightarrow \infty$ [Garrigos and Gower, 2023].

⁷ $c := \min\{\mu_{\text{PL}}/2, \rho/2\}$, $C'_{\text{var}} = C_{\text{var}} \left(L(1 + L/\rho) + \alpha \frac{L^2}{\rho} \right) G^2$ for some numerical constant $C_{\text{var}} > 0$, $C_{\text{dist}} = \sqrt{\left(1 + \frac{L}{\rho}\right)^2 \frac{2}{\mu_{\text{QG}}} + \frac{2}{\alpha\rho}}$, and $C'_{\text{hit}} = 2G \left(\frac{(1+L/\rho)^2}{\sqrt{\mu_{\text{PL}}\mu_{\text{QG}}}} + \frac{1}{\rho} \right) + \frac{G}{L} \left(9 + \frac{L}{\rho} \right)$.

254 be the on-average argument-stability constant from Theorem 3. Then

$$\begin{aligned} \mathbb{E}[R(w_T) - R_n(w_T)] &\leq (1 + L/\rho) G \varepsilon_T, \\ |\mathbb{E}[\Delta^{\text{PD}}(w_T, v_T) - \Delta_n^{\text{PD}}(w_T, v_T)]| &\leq G \varepsilon_T. \end{aligned}$$

255 Combining Corollary 4 and Lemma 5 [Wang et al., 2022], we arrive at Theorem 6, the main result of
 256 this paper, i.e., the *generalization guarantee of Bellman residuals*. In words, the learned Q (i.e., the
 257 corresponding learned w) generalizes: its empirical Bellman residual on the offline dataset closely
 258 matches its expected Bellman residual on the true MDP distribution. This is direct from the fact that
 259 proving $\mathcal{O}(1/n)$ stability for SGDA immediately delivers $\mathcal{O}(1/n)$ generalization bounds for Bellman
 260 residual minimization [Wang et al., 2022].

261 **Theorem 6** (Generalization for the empirical Bellman–residual objective). *Let $(\hat{w}^{(T)}, \hat{v}^{(T)})$ be the*
 262 *parameters returned by ALGORITHM 1 after T SGDA iterations on the empirical objective (5).*
 263 *Define the population and empirical risks*

$$\mathcal{R}(w) := \mathbb{E}_{(s,a) \sim \pi_D, \nu_0} \mathbb{E}_{s' \sim P(\cdot|s,a)} [f(w, \tilde{v}^*; s, a, s')], \quad \hat{\mathcal{R}}_n(w) := \frac{1}{N} \sum_{(s,a,s') \in \mathcal{D}_{\pi_D, \nu_0}} [f(w, \tilde{v}^*; s, a, s')],$$

264 where \tilde{v}^* is the minimizer in (5). Then for ε_T from Theorem 3, we have

$$\begin{aligned} \mathbb{E}[\mathcal{R}(\hat{w}^{(T)}) - \hat{\mathcal{R}}_n(\hat{w}^{(T)})] &\leq (1 + L/\rho) G \varepsilon_T, \\ |\mathbb{E}[\Delta^{\text{PD}}(\hat{w}^{(T)}, \hat{v}^{(T)}) - \Delta_n^{\text{PD}}(\hat{w}^{(T)}, \hat{v}^{(T)})]| &\leq G \varepsilon_T. \end{aligned}$$

265 The generalization guarantee in Theorem 6 is a critical result, confirming that the empirical risk is a
 266 reliable proxy for the true population risk. However, the ultimate measure of success for a learning
 267 algorithm is its performance on the population distribution relative to the best possible model. This is
 268 quantified by the *population excess risk*, which measures the gap $\mathcal{R}(\hat{w}^{(T)}) - \mathcal{R}(w^*)$.

269 **Theorem 7** (Population excess risk). *Let $(\hat{w}^{(T)}, \hat{v}^{(T)})$ be the SGDA iterate outcome of Algorithm 1*
 270 *after T steps and $w^* \in \arg \min_{w \in \mathcal{W}} \mathcal{R}(w)$ its population minimiser. Then, with ε_T from Theorem 3*
 271 *and d_1, d_2 defined in Lemma 2,*

$$\mathbb{E}[\mathcal{R}(\hat{w}^{(T)}) - \mathcal{R}(w^*)] \leq \underbrace{(1 + L/\rho) G \varepsilon_T}_{\text{stability / generalization}} + \underbrace{\frac{d_1}{d_2 + T}}_{\text{optimization error}}$$

272 4 Conclusion and Limitation

273 We studied the statistical behavior of Bellman Residual Minimization (BRM) in offline RL/IRL
 274 through the lens of stability. Exploiting the PL–strongly-concave geometry of the bi-conjugate
 275 formulation, we coupled two SGDA trajectories on neighboring datasets with a single Lyapunov
 276 potential and a ghost-index decoupling device. This yielded an *on-average argument-stability* bound
 277 with $\mathcal{O}(1/n)$ rate (Theorem 3), which directly implies $\mathcal{O}(1/n)$ generalization for BRM (Theorem 6)
 278 and a population excess-risk bound that cleanly decomposes optimization and estimation errors (The-
 279 orem 7). The analysis is constructive, tracks explicit constants in $(L, \rho, \mu_{\text{PL}}, \mu_{\text{QG}}, G)$, accommodates
 280 minibatching, and requires neither variance reduction nor independence assumptions on the sampling
 281 indices. Together with the global convergence of SGDA in parameter space (Lemma 2), these results
 282 close the statistical gap for BRM and improve the sample-complexity exponent over the $\mathcal{O}(n^{-1/2})$
 283 rates known for convex–concave saddle problems.

284 About the future work, this work focuses on finite-sample stability and generalization guarantees
 285 for smooth, entropy-regularized BRM, and therefore does not address genuinely non-smooth BRM
 286 variants or the conversion from Bellman-residual guarantees to policy-performance guarantees
 287 without additional coverage assumptions. We discuss these topics further in Appendix C.6, and
 288 view them as natural next steps toward a more complete statistical theory of BRM-based offline
 289 reinforcement learning.

290 References

- 291 Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-
292 parameterization. In *Proceedings of the 36th International Conference on Machine Learning*,
293 pages 242–252. PMLR, 2019.
- 294 András Antos, Csaba Szepesvári, and Rémi Munos. Learning near-optimal policies with bellman-
295 residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 71
296 (1):89–129, 2008.
- 297 Sanjeev Arora, Simon S. Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of
298 optimization and generalization for overparameterized two-layer neural networks. In *Proceedings*
299 *of the 36th International Conference on Machine Learning*, pages 322–332. PMLR, 2019.
- 300 Leemon Baird. Residual algorithms: Reinforcement learning with function approximation. In
301 *Machine Learning Proceedings 1995*, pages 30–37. Elsevier, 1995.
- 302 Olivier Bousquet and André Elisseeff. Stability and generalization. *Journal of machine learning*
303 *research*, 2(Mar):499–526, 2002.
- 304 Olivier Bousquet, Yegor Klochkov, and Nikita Zhivotovskiy. Sharper bounds for uniformly stable
305 algorithms. In *Proceedings of Thirty Third Conference on Learning Theory*, pages 610–626.
306 PMLR, 2020.
- 307 Zachary Charles and Dimitris Papailiopoulos. Stability and generalization of learning algorithms
308 that converge to global optima. In *International conference on machine learning*, pages 745–754.
309 PMLR, 2018.
- 310 Yixuan Chen, Yubin Shi, Mingzhi Dong, Xiaochen Yang, Dongsheng Li, Yujiang Wang, Robert P.
311 Dick, Qin Lv, Yingying Zhao, Fan Yang, Ning Gu, and Li Shang. Over-parameterized model
312 optimization with polyak-łojasiewicz condition. In *International Conference on Learning Repre-*
313 *sentations*, 2023.
- 314 Bo Dai, Albert Shaw, Lihong Li, Lin Xiao, Niao He, Zhen Liu, Jianshu Chen, and Le Song. Sbeed:
315 Convergent reinforcement learning with nonlinear function approximation. In *International*
316 *conference on machine learning*, pages 1125–1134. PMLR, 2018.
- 317 Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes
318 over-parameterized neural networks. In *International Conference on Learning Representations*,
319 2019.
- 320 Damien Ernst, Pierre Geurts, and Louis Wehenkel. Tree-based batch mode reinforcement learning.
321 *Journal of Machine Learning Research*, 6, 2005.
- 322 Farzan Farnia and Asuman Ozdaglar. Train simultaneously, generalize better: Stability of gradient-
323 based minimax learners. In *Proceedings of the 38th International Conference on Machine Learning*,
324 pages 3174–3185. PMLR, 2021.
- 325 Vitaly Feldman and Jan Vondrak. Generalization bounds for uniformly stable algorithms. *Advances*
326 *in Neural Information Processing Systems*, 31, 2018.
- 327 Guillaume Garrigos and Robert M. Gower. Handbook of convergence theorems for (stochastic)
328 gradient methods. arXiv preprint arXiv:2301.11235, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2301.11235)
329 [2301.11235](https://arxiv.org/abs/2301.11235).
- 330 Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic
331 gradient descent. In *International conference on machine learning*, pages 1225–1234. PMLR,
332 2016.
- 333 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and gener-
334 alization in neural networks. In *Advances in Neural Information Processing Systems*, volume 31,
335 2018.
- 336 Nan Jiang and Tengyang Xie. Offline reinforcement learning in large state spaces: Algorithms and
337 guarantees. *Statistical Science*, 2024.

- 338 Enoch H. Kang, Hema Yoganarasimhan, and Lalit Jain. An empirical risk minimization approach
339 for offline inverse rl and dynamic discrete choice model, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2502.14131)
340 [2502.14131](https://arxiv.org/abs/2502.14131).
- 341 Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-
342 gradient methods under the polyak-lojasiewicz condition. *arXiv preprint arXiv:1608.04636*, 2016.
343 URL <https://arxiv.org/abs/1608.04636>.
- 344 Ilja Kuzborskij and Christoph H. Lampert. Data-dependent stability of stochastic gradient descent. In
345 *Proceedings of the 35th International Conference on Machine Learning*, pages 2815–2824. PMLR,
346 2018.
- 347 Jaehoon Lee, Lechao Xiao, Samuel S. Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-
348 Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models
349 under gradient descent. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- 350 Yunwen Lei and Yiming Ying. Fine-grained analysis of stability and generalization for stochastic
351 gradient descent. *Proceedings of the 37th International Conference on Machine Learning*, pages
352 5809–5819, 2020.
- 353 Yunwen Lei, Zhenhuan Yang, Tianbao Yang, and Yiming Ying. Stability and generalization of
354 stochastic gradient methods for minimax problems. In Marina Meilă and Tong Zhang, edi-
355 tors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of
356 *Proceedings of Machine Learning Research*, pages 6175–6186. PMLR, 2021. URL [https://](https://proceedings.mlr.press/v139/lei21b.html)
357 proceedings.mlr.press/v139/lei21b.html.
- 358 Chaoyue Liu, Libin Zhu, and Mikhail Belkin. Loss landscapes and optimization in over-parameterized
359 non-linear systems and neural networks. *Applied and Computational Harmonic Analysis*, 59:
360 85–116, 2022.
- 361 Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon: Infinite-
362 horizon off-policy estimation. In *Advances in Neural Information Processing Systems*, 2018.
- 363 Ofir Nachum, Yinlam Chow, Bo Dai, and Lihong Li. Dualdice: Behavior-agnostic estimation of
364 discounted stationary distribution corrections. In *Advances in Neural Information Processing*
365 *Systems*, volume 32, 2019.
- 366 Yurii Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Springer, 2004.
- 367 Andrew Patterson, Adam White, and Martha White. A generalized projected bellman error for
368 off-policy value estimation in reinforcement learning. *Journal of Machine Learning Research*, 23
369 (145):1–61, 2022.
- 370 Boris T. Polyak. Gradient methods for minimizing functionals. *USSR Computational Mathematics*
371 *and Mathematical Physics*, 3(4):864–878, 1963.
- 372 John Rust. Structural estimation of markov decision processes. *Handbook of econometrics*, 4:
373 3081–3143, 1994.
- 374 Mahdi Soltanolkotabi, Adel Javanmard, and Jason D. Lee. Theoretical insights into the optimization
375 landscape of over-parameterized shallow neural networks. *IEEE Transactions on Information*
376 *Theory*, 65(2):742–769, 2018.
- 377 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2
378 edition, 2018.
- 379 John N Tsitsiklis and Benjamin Van Roy. Feature-based methods for large scale dynamic program-
380 ming. *Machine Learning*, 22(1):59–94, 1996.
- 381 Hado Van Hasselt, Yotam Doron, Florian Strub, Matteo Hessel, Nicolas Sonnerat, and Joseph
382 Modayil. Deep reinforcement learning and the deadly triad. *arXiv preprint arXiv:1812.02648*,
383 2018.

- 384 Puyu Wang, Yunwen Lei, Yiming Ying, and Ding-Xuan Zhou. Stability and generalization for
385 markov chain stochastic gradient methods. In *Advances in Neural Information Processing Systems*,
386 volume 35, 2022. URL [https://proceedings.neurips.cc/paper_files/paper/2022/
387 hash/f61538f83b0f19f9306d9d801c15f41c-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2022/hash/f61538f83b0f19f9306d9d801c15f41c-Abstract-Conference.html).
- 388 Ziqing Xu, Hancheng Min, Salma Tarmoun, Enrique Mallada, and René Vidal. A local polyak-
389 łojasiewicz and descent lemma of gradient descent for overparametrized linear models. *Transac-
390 tions on Machine Learning Research*, 2025.
- 391 Junchi Yang, Negar Kiyavash, and Niao He. Global convergence and variance-reduced optimization
392 for a class of nonconvex-nonconcave minimax problems, 2020. URL [https://arxiv.org/abs/
393 2002.09621](https://arxiv.org/abs/2002.09621).
- 394 Ruiyi Zhang, Bo Dai, Lihong Li, and Dale Schuurmans. Gendice: Generalized offline estimation of
395 stationary values. In *International Conference on Learning Representations*, 2020.
- 396 Siqi Zhang, Yifan Hu, Liang Zhang, and Niao He. Generalization bounds of nonconvex-(strongly)-
397 concave stochastic minimax optimization. In *Proceedings of The 27th International Conference on
398 Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*,
399 pages 694–702. PMLR, 2024.
- 400 Miaoxi Zhu, Li Shen, Bo Du, and Dacheng Tao. Stability and generalization of the decentralized
401 stochastic gradient descent ascent algorithm. In *Advances in Neural Information Processing
402 Systems*, volume 36, 2023.

Symbol	Meaning
\mathcal{S}	State space
\mathcal{A}	Action space
$P(\cdot s, a)$	Transition kernel over the next state from the state-action pair (s, a)
$r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$	Deterministic reward function
$\Delta_{\mathcal{S}}^{\mathcal{A}}$	Set of stationary Markov policies
$\pi \in \Delta_{\mathcal{S}}^{\mathcal{A}}$	Policy
$\mathcal{H}(q)$	Shannon entropy of an action distribution q
β	Discount factor
η	Entropy regularization coefficient
\mathcal{T}	Bellman optimality operator
$\hat{\mathcal{T}}$	Sampled Bellman operator
$\delta_Q(s, a, s')$	Temporal-Difference error
Q^*	(The unique) solution to Bellman optimality equation
$z_i = (s_i, a_i, s'_i)$	i -th sample
$\mathcal{D} = \{z_i\}_{i=1}^n$	Dataset
$f(w, v; z_i)$	Per-sample objective
$\mathcal{D}^{(i)}$	neighbouring dataset
g_t^w, g_t^v	per-sample gradients.
η_t	SGDA learning rate at time t
$F_D(w, v)$	$\frac{1}{n} \sum_{j=1}^n f(w, v; z_j^{\mathcal{D}})$
$\Phi_D(w)$	$\max_v F_D(w, v)$
Φ_D^*	$\min_w \Phi_D(w)$
L	Smoothness (Condition (C1))
μ_{PL}	Polyak–Łojasiewicz condition constant (Condition (C2))
μ_{QG}	Quadratic Growth constant (Condition (C3))
ρ	Strong concavity constant (Condition (C4))
G	Per-sample gradient bound (Condition (C5))
$\Psi_{\alpha, \mathcal{D}}(w, v)$	Lyapunov potential, $\Phi_D(w) - \Phi_D^* + \alpha \cdot (\Phi_D(w) - F_D(w, v))$
$A(w)$	$\Phi_D(w) - \Phi_D^*$
$\Gamma(w, v)$	$\Phi_D(w) - F_D(w, v)$
$R(w), R_n(w)$	Primal Risk (Definition 4)
$\Delta^{PD}, \Delta_n^{PD}$	Weak primal-dual risk (Definition 5)

Table 1: Notations

405 B Technical Proofs

406 B.1 Supporting Lemmas for proving Theorem 3

407 **Lemma 8** (Mismatch control). *Let $v_D^*(w) := \arg \max_v F_D(w, v)$. Under (C1) and (C4),*

$$\Delta(w, v) := \nabla_w F_D(w, v) - \nabla \Phi_D(w) = \nabla_w F_D(w, v) - \nabla_w F_D(w, v_D^*(w))$$

408 *satisfies $\|\Delta(w, v)\| \leq L \|v - v_D^*(w)\|$. Moreover,*

$$\|v - v_D^*(w)\|^2 \leq \frac{2}{\rho} (\Phi_D(w) - F_D(w, v)).$$

409 *Proof.* L -Lipschitzness of $\nabla_w F_D(w, \cdot)$ yields the first bound. For $g(\cdot) := F_D(w, \cdot)$, ρ -strong
410 concavity implies $g(v^*) - g(v) \geq (\rho/2) \|v - v^*\|^2$ [Nesterov, 2004], which gives the second inequality.411 \square

412 **Lemma 9** (Smoothness of the value function Φ). *Assume (C1) (joint L -smoothness) and (C4)*
 413 *(ρ -strong concavity in v). Then, for each dataset D :*

414 (i) *the maximizer $v_D^*(w) := \arg \max_v F_D(w, v)$ is well-defined and (L/ρ) -Lipschitz:*

$$\|v_D^*(w) - v_D^*(u)\| \leq \frac{L}{\rho} \|w - u\| \quad \forall w, u,$$

415 (ii) $\Phi_D(w) := \max_v F_D(w, v)$ *is differentiable with $\nabla \Phi_D(w) = \nabla_w F_D(w, v_D^*(w))$ (Dan-*
 416 *skin),*

417 (iii) *and Φ_D is L_Φ -smooth with*

$$\|\nabla \Phi_D(w) - \nabla \Phi_D(u)\| \leq L \left(1 + \frac{L}{\rho}\right) \|w - u\| \quad \forall w, u.$$

418 *Proof. (i) Lipschitzness of $v_D^*(\cdot)$.* By (C4), for each fixed w the map $v \mapsto F_D(w, v)$ is ρ -
 419 strongly concave, so $v_D^*(w)$ is unique. Using the first-order conditions $\nabla_v F_D(w, v_D^*(w)) = 0$
 420 and $\nabla_v F_D(u, v_D^*(u)) = 0$, write

$$\begin{aligned} 421 \quad 0 &= \nabla_v F_D(w, v_D^*(w)) - \nabla_v F_D(u, v_D^*(u)) \\ &= \underbrace{[\nabla_v F_D(w, v_D^*(w)) - \nabla_v F_D(w, v_D^*(u))]}_{(A)} + \underbrace{[\nabla_v F_D(w, v_D^*(u)) - \nabla_v F_D(u, v_D^*(u))]}_{(B)}. \end{aligned}$$

422 Strong concavity makes $\nabla_v F_D(w, \cdot)$ ρ -strongly monotone, so $\langle (A), v_D^*(w) - v_D^*(u) \rangle \leq -\rho \|v_D^*(w) -$
 423 $v_D^*(u)\|^2$. Joint L -smoothness gives $\|(B)\| \leq L \|w - u\|$. Taking inner products with $v_D^*(w) - v_D^*(u)$
 424 yields

$$\rho \|v_D^*(w) - v_D^*(u)\| \leq \|(B)\| \leq L \|w - u\| \quad \Rightarrow \quad \|v_D^*(w) - v_D^*(u)\| \leq \frac{L}{\rho} \|w - u\|.$$

425 (ii) By joint smoothness, Danskin's theorem applies: $\nabla \Phi_D(w) = \nabla_w F_D(w, v_D^*(w))$.

426 (iii) *Smoothness of Φ_D .* For any w, u ,

$$\begin{aligned} &\|\nabla \Phi_D(w) - \nabla \Phi_D(u)\| \\ &= \|\nabla_w F_D(w, v_D^*(w)) - \nabla_w F_D(u, v_D^*(u))\| \\ &\leq \underbrace{\|\nabla_w F_D(w, v_D^*(w)) - \nabla_w F_D(w, v_D^*(u))\|}_{\leq L \|v_D^*(w) - v_D^*(u)\|} + \underbrace{\|\nabla_w F_D(w, v_D^*(u)) - \nabla_w F_D(u, v_D^*(u))\|}_{\leq L \|w - u\|} \\ &\leq L \left(\frac{L}{\rho} + 1\right) \|w - u\|. \end{aligned}$$

427 The first two inequalities use joint L -smoothness (of $\nabla_w F_D$ in both arguments); the last uses part (i).
 428 Thus $L_\Phi \leq L(1 + L/\rho)$. \square

429 **Lemma 10** (Cross-dataset gradient sensitivity for Φ). *Under (C1), (C4) and (C5), for any w, u and*
 430 *any 1-sample replacement $D \rightarrow D^{(i)}$,*

$$\|\nabla \Phi_D(w) - \nabla \Phi_{D^{(i)}}(u)\| \leq L_\Phi \|w - u\| + \frac{2G}{n} \left(1 + \frac{L}{\rho}\right).$$

431 *Proof.* By Lemma 9, $v_D^*(\cdot)$ is (L/ρ) -Lipschitz, Φ_D is L_Φ -smooth, and $\nabla \Phi_D(w) =$
 432 $\nabla_w F_D(w, v_D^*(w))$. Hence

$$\|\nabla \Phi_D(w) - \nabla \Phi_{D^{(i)}}(u)\| \leq \underbrace{\|\nabla \Phi_D(w) - \nabla \Phi_D(u)\|}_{\leq L_\Phi \|w - u\|} + \underbrace{\|\nabla \Phi_D(u) - \nabla \Phi_{D^{(i)}}(u)\|}_{(A)}.$$

433 We bound (A) by inserting and subtracting $v_{D^{(i)}}^*(u)$ and using joint L -smoothness of F :

$$\begin{aligned} (A) &= \|\nabla_w F_D(u, v_D^*(u)) - \nabla_w F_{D^{(i)}}(u, v_{D^{(i)}}^*(u))\| \\ &\leq \underbrace{\|\nabla_w F_D(u, v_D^*(u)) - \nabla_w F_D(u, v_{D^{(i)}}^*(u))\|}_{\leq L \|v_D^*(u) - v_{D^{(i)}}^*(u)\|} + \underbrace{\|\nabla_w F_D(u, v_{D^{(i)}}^*(u)) - \nabla_w F_{D^{(i)}}(u, v_{D^{(i)}}^*(u))\|}_{\leq \frac{2G}{n} \text{ by Condition (C5)}}, \end{aligned}$$

434 so that

$$(A) \leq L \|v_D^*(u) - v_{D^{(i)}}^*(u)\| + \frac{2G}{n}.$$

435 Since $v \mapsto F_D(u, v)$ is ρ -strongly concave, the gradient map is ρ -strongly monotone, which gives
436 the error bound with the normal cone $N_{\mathcal{V}}(\cdot)$:

$$\|v - v_D^*(u)\| \leq \frac{1}{\rho} \text{dist}(\nabla_v F_D(u, v), N_{\mathcal{V}}(v)) \quad \forall v \in \mathcal{V}.$$

437 Apply this at $v = v_{D^{(i)}}^*(u)$. There are two cases.

438 (i) *Unconstrained (or interior) maximizer.* Then $N_{\mathcal{V}}(v_{D^{(i)}}^*(u)) = \{0\}$, so

$$\|v_D^*(u) - v_{D^{(i)}}^*(u)\| \leq \frac{1}{\rho} \|\nabla_v F_D(u, v_{D^{(i)}}^*(u))\| = \frac{1}{\rho} \|\nabla_v F_D(u, v_{D^{(i)}}^*(u)) - \nabla_v F_{D^{(i)}}(u, v_{D^{(i)}}^*(u))\|.$$

439 (ii) *Constrained boundary maximizer.* By KKT optimality,

$$0 \in -\nabla_v F_{D^{(i)}}(u, v_{D^{(i)}}^*(u)) + N_{\mathcal{V}}(v_{D^{(i)}}^*(u)) \iff \nabla_v F_{D^{(i)}}(u, v_{D^{(i)}}^*(u)) \in N_{\mathcal{V}}(v_{D^{(i)}}^*(u)).$$

440 Hence, choosing $\xi := \nabla_v F_{D^{(i)}}(u, v_{D^{(i)}}^*(u)) \in N_{\mathcal{V}}(v_{D^{(i)}}^*(u))$, we obtain

$$\text{dist}(\nabla_v F_D(u, v_{D^{(i)}}^*(u)), N_{\mathcal{V}}(v_{D^{(i)}}^*(u))) \leq \|\nabla_v F_D(u, v_{D^{(i)}}^*(u)) - \nabla_v F_{D^{(i)}}(u, v_{D^{(i)}}^*(u))\|.$$

441 In either case, using the single-sample replacement identity

$$\nabla_v F_D(u, v) - \nabla_v F_{D^{(i)}}(u, v) = \frac{1}{n} \left(\nabla_v f(u, v; z_i) - \nabla_v f(u, v; \tilde{z}_i) \right),$$

442 the triangle inequality and Condition (C5) give

$$\begin{aligned} \|v_D^*(u) - v_{D^{(i)}}^*(u)\| &\leq \frac{1}{\rho} \text{dist}(\nabla_v F_D(u, v_{D^{(i)}}^*(u)), N_{\mathcal{V}}(v_{D^{(i)}}^*(u))) \\ &\leq \frac{1}{\rho} \|\nabla_v F_D(u, v_{D^{(i)}}^*(u)) - \nabla_v F_{D^{(i)}}(u, v_{D^{(i)}}^*(u))\| \\ &= \frac{1}{\rho} \cdot \frac{1}{n} \|\nabla_v f(u, v_{D^{(i)}}^*(u); z_i) - \nabla_v f(u, v_{D^{(i)}}^*(u); \tilde{z}_i)\| \\ &\leq \frac{1}{\rho} \cdot \frac{1}{n} \left(\|\nabla_v f(u, v_{D^{(i)}}^*(u); z_i)\| + \|\nabla_v f(u, v_{D^{(i)}}^*(u); \tilde{z}_i)\| \right) \\ &\leq \frac{2G}{\rho n}. \end{aligned}$$

443 Plugging this into the bound for (A) yields

$$(A) \leq L \cdot \frac{2G}{\rho n} + \frac{2G}{n} = \frac{2G}{n} \left(1 + \frac{L}{\rho} \right),$$

444 and therefore

$$\|\nabla \Phi_D(w) - \nabla \Phi_{D^{(i)}}(u)\| \leq L_{\Phi} \|w - u\| + \frac{2G}{n} \left(1 + \frac{L}{\rho} \right).$$

445 □

446 **Lemma 11** (Deterministic potential \Rightarrow distance). *Under (C3) and (C4), for any $\alpha > 0$ and any*
447 *(w, v),*

$$\|w - x_D^*\| + \|v - v_D^*\| \leq \sqrt{\left(1 + \frac{L}{\rho}\right)^2 \frac{2}{\mu_{\text{QG}}} + \frac{2}{\alpha \rho}} \cdot \sqrt{\Psi_{\alpha, D}(w, v)}.$$

448 *Proof.* Using $\|v - v_D^*\| \leq \|v - v_D^*(w)\| + \|v_D^*(w) - v_D^*\|$ and the (L/ρ) -Lipschitzness of $w \mapsto v_D^*(w)$,

$$\|w - x_D^*\| + \|v - v_D^*\| \leq \left(1 + \frac{L}{\rho}\right) \|w - x_D^*\| + \|v - v_D^*(w)\|.$$

$$\|w - x_D^*\| \leq \sqrt{\frac{2}{\mu_{\text{QG}}}(\Phi_D(w) - \Phi_D^*)}, \quad \|v - v_D^*(w)\| \leq \sqrt{\frac{2}{\rho}(\Phi_D(w) - F_D(w, v))},$$

450 where $v_D^*(w) \in \arg \max_u F_D(w, u)$ and $\Phi_D(w) = \max_u F_D(w, u)$.

451 Using the weighted Cauchy–Schwarz inequality with any $\alpha > 0$,

$$\begin{aligned} & (1 + \frac{L}{\rho})\sqrt{\frac{2}{\mu_{\text{QG}}}(\Phi_D(w) - \Phi_D^*)} + \sqrt{\frac{2}{\rho}(\Phi_D(w) - F_D(w, v))} \\ & \leq \sqrt{(\Phi_D(w) - \Phi_D^*) + \alpha(\Phi_D(w) - F_D(w, v))} \sqrt{\frac{2(1+L/\rho)^2}{\mu_{\text{QG}}} + \frac{2}{\alpha\rho}}. \end{aligned}$$

452 Noting that $\Psi_{\alpha, D}(w, v) := (\Phi_D(w) - \Phi_D^*) + \alpha(\Phi_D(w) - F_D(w, v))$, we obtain

$$(1 + \frac{L}{\rho})\sqrt{\frac{2}{\mu_{\text{QG}}}(\Phi_D(w) - \Phi_D^*)} + \sqrt{\frac{2}{\rho}(\Phi_D(w) - F_D(w, v))} \leq \sqrt{\Psi_{\alpha, D}(w, v)} \sqrt{\frac{2(1+L/\rho)^2}{\mu_{\text{QG}}} + \frac{2}{\alpha\rho}}.$$

453 □

454 **Lemma 12** (Saddle-point sensitivity). *Let $D^{(i)}$ be obtained from D by replacing one sample. Assume*
 455 *(C1) (joint L -smoothness), (C4) (ρ -strong concavity in v), (C5) (per-sample gradients bounded by G),*
 456 *(C2)–(C3) (PL and QG for Φ_D), and (S3) (the same constants hold for D and $D^{(i)}$). Let (x_D^*, v_D^*)*
 457 *and $(x_{D^{(i)}}^*, v_{D^{(i)}}^*)$ be the selected saddle points from Condition (C6) for D and $D^{(i)}$, respectively.*
 458 *Then Then*

$$\|x_D^* - x_{D^{(i)}}^*\| + \|v_D^* - v_{D^{(i)}}^*\| \leq \frac{2G}{n} \left(\frac{(1+L/\rho)^2}{\sqrt{\mu_{\text{PL}}\mu_{\text{QG}}}} + \frac{1}{\rho} \right).$$

459 *Proof.* At $(x_{D^{(i)}}^*, v_{D^{(i)}}^*)$, only one summand in F_D differs from $F_{D^{(i)}}$, and per-sample gradients are
 460 bounded by G . Hence,

$$\begin{aligned} \|\nabla_w F_D(x_{D^{(i)}}^*, v_{D^{(i)}}^*)\| &= \|\nabla_w F_D(x_{D^{(i)}}^*, v_{D^{(i)}}^*) - \nabla_w F_D^{(i)}(x_{D^{(i)}}^*, v_{D^{(i)}}^*)\| \leq \frac{2G}{n} \\ \|\nabla_v F_D(x_{D^{(i)}}^*, v_{D^{(i)}}^*)\| &= \|\nabla_v F_D(x_{D^{(i)}}^*, v_{D^{(i)}}^*) - \nabla_v F_D^{(i)}(x_{D^{(i)}}^*, v_{D^{(i)}}^*)\| \leq \frac{2G}{n}. \end{aligned}$$

461 Using Lemma 8,

$$\|\nabla \Phi_D(x_{D^{(i)}}^*)\| \leq \|\nabla_w F_D(x_{D^{(i)}}^*, v_{D^{(i)}}^*)\| + L \|v_{D^{(i)}}^* - v_D^*(x_{D^{(i)}}^*)\| \leq \frac{2G}{n}(1+L/\rho),$$

462 where the last step uses ρ -strong concavity to get $\|v_{D^{(i)}}^* - v_D^*(x_{D^{(i)}}^*)\| \leq \frac{1}{\rho} \|\nabla_v F_D(x_{D^{(i)}}^*, v_{D^{(i)}}^*)\|$.

463 Since Φ_D satisfies PL and QG,

$$\|x_D^* - x_{D^{(i)}}^*\| \leq \frac{1}{\sqrt{\mu_{\text{PL}}\mu_{\text{QG}}}} \|\nabla \Phi_D(x_{D^{(i)}}^*)\| \leq \frac{2G}{n} \cdot \frac{1+L/\rho}{\sqrt{\mu_{\text{PL}}\mu_{\text{QG}}}}.$$

464 For the dual variable, split and use the (L/ρ) -Lipschitzness of $w \mapsto v_D^*(w)$ (Lemma 9) and strong
 465 concavity:

$$\begin{aligned} \|v_D^* - v_{D^{(i)}}^*\| &\leq \|v_D^*(x_D^*) - v_{D^{(i)}}^*(x_{D^{(i)}}^*)\| + \|v_{D^{(i)}}^*(x_{D^{(i)}}^*) - v_{D^{(i)}}^*(x_{D^{(i)}}^*)\| \\ &\leq \frac{L}{\rho} \|x_D^* - x_{D^{(i)}}^*\| + \frac{1}{\rho} \|\nabla_v F_D(x_{D^{(i)}}^*, v_{D^{(i)}}^*)\| \\ &\leq \frac{2G}{n} \left(\frac{L}{\rho} \cdot \frac{1+L/\rho}{\sqrt{\mu_{\text{PL}}\mu_{\text{QG}}}} + \frac{1}{\rho} \right), \end{aligned}$$

466 where the last inequality uses the same single-sample replacement and gradient bound as above (cf.
 467 Lemma 10). Adding the two bounds yields the claim. □

468 **Lemma 13** (Coupled one-step bounds). *Assume (C1)–(S3). Define*

$$\Xi_t := \Psi_{\alpha, D}(w_t, v_t) + \Psi_{\alpha, D^{(i)}}(w'_t, v'_t), \quad D_t := \|w_t - w'_t\| + \|v_t - v'_t\|.$$

469 Let $c := \min\{\mu_{\text{PL}}/2, \rho/2\}$ and $L_\Phi \leq L(1 + L/\rho)$ (Lemma 9). Then, for $\tilde{B}_0 = \tilde{B}_1 = 8$ and a
 470 constant $C_{\text{var}} > 0$ that depends only on L, ρ , we have

$$(P) \quad \mathbb{E}[\Xi_{t+1} \mid \mathcal{F}_t] \leq (1 - c\eta_t)\Xi_t + A\eta_t^2 \\ + C_{\text{leak}}\eta_t D_t + C_* \frac{G}{n}\eta_t + C_{\text{lin}}G^2\eta_t + \eta_t(\mathbf{1}\{i_t = i\} + \mathbf{1}\{\hat{i}_t = i\})\tilde{B}_0G^2, \quad (6)$$

$$(D_{\text{weak}}) \quad \mathbb{E}[D_{t+1} \mid \mathcal{F}_t] \leq (1 + \kappa\eta_t)D_t + a\eta_t\sqrt{\Xi_t} + \eta_t\mathbf{1}\{i_t = i\}\tilde{B}_1G + \eta_t\frac{2G}{n}\left(1 + \frac{L}{\rho}\right), \quad (7)$$

471 where

$$A := C_{\text{var}}\left(L(1 + L/\rho) + \alpha L^2/\rho\right)G^2, \quad C_{\text{leak}} := 2(2L + L_\Phi),$$

$$a := 2C_{\text{dist}}(L + L_\Phi), \quad \kappa := 2L, \quad C_* := 4(1 + L/\rho), \quad C_{\text{dist}} := \sqrt{\left(1 + \frac{L}{\rho}\right)^2 \frac{2}{\mu_{\text{QG}}} + \frac{2}{\alpha\rho}},$$

$$\text{and} \quad C_{\text{lin}} := 12 + \frac{10L^2}{\rho^2}.$$

472 **Why introduce a ghost index?** Writing and conditioning on the past, the primal descent term
 473 $-\eta_t\langle\nabla\Phi(w_t), \nabla_w f(w_t, v_t; z_{i_t})\rangle$ has no sign we can control, since i_t is already revealed in \mathcal{F}_t . We
 474 therefore add–subtract a ghost gradient and take $\mathbb{E}_{\hat{i}_t}[\cdot \mid \mathcal{F}_t]$:

$$-\eta_t\langle\nabla\Phi(w_t), g_t^w\rangle = -\eta_t\langle\nabla\Phi(w_t), \hat{g}_t^w\rangle - \eta_t\langle\nabla\Phi(w_t), g_t^w - \hat{g}_t^w\rangle.$$

475 The first term becomes the correct drift $-\eta_t\langle\nabla\Phi(w_t), \nabla_w F(w_t, v_t)\rangle$, which contracts by PL for Φ ,
 476 while the second is a centered correction that we bound by Young’s inequality and the gradient bound
 477 $\|g_t^w - \hat{g}_t^w\| \leq 2G$, yielding an $O(\eta_t G^2)$ remainder absorbed into the variance term. An identical
 478 maneuver applies to the dual-gap part Γ . This device avoids any i.i.d. sampling assumption and yields
 479 the one-step recursion (P) with a contractive factor and small, explicit noise.

480 *Proof.* All expectations are conditional on $\mathcal{F}_t := \sigma((w_s, v_s, w'_s, v'_s, i_s)_{0 \leq s \leq t})$. Write $\gamma := \eta_t$. Fix
 481 $D \in \{\mathcal{D}, \mathcal{D}^{(i)}\}$ and abbreviate $F := F_D$, $\Phi := \Phi_D$, $\Psi_\alpha := \Psi_{\alpha, D}$, $\Gamma(w, v) := \Phi(w) - F(w, v)$.
 482 Set $g^w := \nabla_w f(w_t, v_t; z_{i_t})$, $g^v := \nabla_v f(w_t, v_t; z_{i_t})$,

$$w^+ = w_t - \gamma g^w, \quad v^+ = v_t + \gamma g^v,$$

483 and introduce a ghost index \hat{i}_t independent of \mathcal{F}_t with $\hat{g}^w := \nabla_w f(w_t, v_t; z_{\hat{i}_t})$, $\hat{g}^v :=$
 484 $\nabla_v f(w_t, v_t; z_{\hat{i}_t})$.

485 **(a) Primal part Φ .** Recall $w^+ = w_t - \gamma g^w$ with $g^w := \nabla_w f(w_t, v_t; z_{i_t})$ and let $\gamma := \eta_t$. By
 486 L_Φ -smoothness of Φ ,

$$\Phi(w^+) - \Phi(w_t) \leq -\gamma\langle\nabla\Phi(w_t), g^w\rangle + \frac{L_\Phi}{2}\gamma^2\|g^w\|^2.$$

487 Insert and subtract a ghost gradient $\hat{g}^w := \nabla_w f(w_t, v_t; z_{\hat{i}_t})$ with $\hat{i}_t \perp \mathcal{F}_t$, then take $\mathbb{E}_{\hat{i}_t}[\cdot \mid \mathcal{F}_t]$. Using
 488 $\mathbb{E}_{\hat{i}_t}[\hat{g}^w \mid \mathcal{F}_t] = \nabla_w F(w_t, v_t)$ and $\|g^w\| \leq G$,

$$\mathbb{E}_{\hat{i}_t}[\Phi(w^+) - \Phi(w_t) \mid \mathcal{F}_t] \leq -\gamma\langle\nabla\Phi(w_t), \nabla_w F(w_t, v_t)\rangle - \gamma\langle\nabla\Phi(w_t), g^w - \nabla_w F(w_t, v_t)\rangle + \frac{L_\Phi}{2}\gamma^2G^2. \quad (8)$$

489 For the centered correction, Young’s inequality with $\|g^w - \nabla_w F(w_t, v_t)\| \leq 2G$ yields

$$\gamma|\langle\nabla\Phi(w_t), g^w - \nabla_w F(w_t, v_t)\rangle| \leq \frac{1}{4}\gamma\|\nabla\Phi(w_t)\|^2 + 4\gamma G^2. \quad (9)$$

490 For the drift, write $\Delta_t := \nabla_w F(w_t, v_t) - \nabla\Phi(w_t)$. Then

$$-\langle\nabla\Phi(w_t), \nabla_w F(w_t, v_t)\rangle = -\|\nabla\Phi(w_t)\|^2 - \langle\nabla\Phi(w_t), \Delta_t\rangle \leq -\frac{1}{2}\|\nabla\Phi(w_t)\|^2 + \frac{1}{2}\|\Delta_t\|^2. \quad (10)$$

491 By Lemma 8 and ρ -strong concavity in v ,

$$\|\Delta_t\| \leq L \|v_t - v^*(w_t)\| \leq L \sqrt{\frac{2}{\rho} \Gamma(w_t, v_t)} \quad \Rightarrow \quad \frac{1}{2} \|\Delta_t\|^2 \leq \frac{L^2}{\rho} \Gamma(w_t, v_t), \quad (11)$$

492 where $\Gamma(w, v) := \Phi(w) - F(w, v)$. Combining (8), (9), (10), and (11),

$$\mathbb{E}_{z_t}[\Phi(w^+) - \Phi(w_t) \mid \mathcal{F}_t] \leq -\frac{1}{4} \gamma \|\nabla \Phi(w_t)\|^2 + \frac{L^2}{\rho} \gamma \Gamma(w_t, v_t) + \frac{L_{\Phi}}{2} \gamma^2 G^2 + 4 \gamma G^2.$$

493 Finally, by the PL inequality $\frac{1}{2} \|\nabla \Phi(w_t)\|^2 \geq \mu_{\text{PL}}(\Phi(w_t) - \Phi^*)$,

$$-\frac{1}{4} \gamma \|\nabla \Phi(w_t)\|^2 \leq -\frac{\mu_{\text{PL}}}{2} \gamma (\Phi(w_t) - \Phi^*),$$

494 hence

$$\mathbb{E}_{z_t}[\Phi(w^+) - \Phi(w_t) \mid \mathcal{F}_t] \leq -\frac{\mu_{\text{PL}}}{2} \gamma (\Phi(w_t) - \Phi^*) + \frac{L^2}{\rho} \gamma \Gamma(w_t, v_t) + \frac{L_{\Phi}}{2} \gamma^2 G^2 + 4 \gamma G^2.$$

495 Since $\Phi(w_t) - \Phi^*$ is \mathcal{F}_t -measurable, adding it to both sides yields

$$\mathbb{E}_{z_t}[\Phi(w^+) - \Phi^* \mid \mathcal{F}_t] \leq (1 - \frac{\mu_{\text{PL}}}{2} \gamma) (\Phi(w_t) - \Phi^*) + \frac{L^2}{\rho} \gamma \Gamma(w_t, v_t) + \frac{L_{\Phi}}{2} \gamma^2 G^2 + 4 \gamma G^2.$$

496 **(b) Dual-gap part Γ .** Recall $\Gamma(w, v) := \Phi_D(w) - F_D(w, v)$ and write $v^+ = v_t + \gamma g_t^v$ with
497 $g_t^v := \nabla_v f(w_t, v_t; z_{i_t})$. By Lemma 14, for fixed w_t ,

$$\Gamma(w_t, v^+) \leq (1 - 2\rho\gamma + \rho L \gamma^2) \Gamma(w_t, v_t). \quad (12)$$

498 Let $w^+ := w_t - \gamma g_t^w$ with $g_t^w := \nabla_w f(w_t, v_t; z_{i_t})$. By L_{Φ} -smoothness of Φ_D and joint L -
499 smoothness of F_D ,

$$\Gamma(w^+, v^+) \leq \Gamma(w_t, v^+) - \gamma \langle \nabla \Phi_D(w_t), g_t^w \rangle + \gamma \langle \nabla_w F_D(w_t, v^+), g_t^w \rangle + \frac{L_{\Phi} + L}{2} \gamma^2 G^2.$$

500 Insert and subtract a ghost gradient $\hat{g}_t^w := \nabla_w f(w_t, v_t; z_{i_t}^{\dagger})$, take $\mathbb{E}_{z_t}[\cdot \mid \mathcal{F}_t]$, and use $\mathbb{E}_{z_t}[\hat{g}_t^w \mid \mathcal{F}_t] =$
501 $\nabla_w F_D(w_t, v_t)$. Then

$$\begin{aligned} \mathbb{E}_{z_t}[\Gamma(w^+, v^+) \mid \mathcal{F}_t] &\leq \Gamma(w_t, v^+) - \gamma \langle \nabla \Phi_D(w_t), \nabla_w F_D(w_t, v_t) \rangle \\ &\quad + \gamma \langle \nabla_w F_D(w_t, v^+) - \nabla \Phi_D(w_t), \nabla_w F_D(w_t, v_t) \rangle \\ &\quad + \gamma \langle \nabla_w F_D(w_t, v^+) - \nabla \Phi_D(w_t), g_t^w - \nabla_w F_D(w_t, v_t) \rangle \\ &\quad + \frac{L_{\Phi} + L}{2} \gamma^2 G^2. \end{aligned}$$

502 Decompose $\nabla_w F_D(w_t, v^+) - \nabla \Phi_D(w_t) = \Delta_t + (\nabla_w F_D(w_t, v^+) - \nabla_w F_D(w_t, v_t))$, where $\Delta_t :=$
503 $\nabla_w F_D(w_t, v_t) - \nabla \Phi_D(w_t)$. Using Lemma 8, $\|\Delta_t\| \leq L \sqrt{2\Gamma(w_t, v_t)}/\rho$, joint L -smoothness,
504 $\|\nabla_w F_D(w_t, v^+) - \nabla_w F_D(w_t, v_t)\| \leq L\gamma G$, and $\|g_t^w - \nabla_w F_D(w_t, v_t)\| \leq 2G$,

$$\begin{aligned} \gamma |\langle \Delta_t, g_t^w - \nabla_w F_D(w_t, v_t) \rangle| &\leq \frac{\rho}{4} \gamma \Gamma(w_t, v_t) + \frac{8L^2}{\rho^2} \gamma G^2, \\ \gamma |\langle \nabla_w F_D(w_t, v^+) - \nabla_w F_D(w_t, v_t), g_t^w - \nabla_w F_D(w_t, v_t) \rangle| &\leq 2L \gamma^2 G^2, \\ \gamma |\langle \Delta_t, \nabla_w F_D(w_t, v_t) \rangle| &\leq \frac{\rho}{4} \gamma \Gamma(w_t, v_t) + \frac{2L^2}{\rho^2} \gamma G^2, \\ \gamma |\langle \nabla_w F_D(w_t, v^+) - \nabla_w F_D(w_t, v_t), \nabla_w F_D(w_t, v_t) \rangle| &\leq L \gamma^2 G^2. \end{aligned}$$

505 Combining the bounds with the ascent contraction (12), we can group the contributions as follows:

506 (i) From the primal descent step we obtain

$$-\frac{1}{4} \gamma \|\nabla \Phi_D(w_t)\|^2 \leq -\frac{\mu_{\text{PL}}}{2} \gamma (\Phi_D(w_t) - \Phi_D^*),$$

507 where the inequality uses the PL condition. From the dual ascent we obtain

$$(1 - 2\rho\gamma + \rho L \gamma^2) \Gamma(w_t, v_t) = \Gamma(w_t, v_t) - 2\rho\gamma \Gamma(w_t, v_t) + O(\gamma^2),$$

508 so the leading negative part is $-2\rho\gamma \Gamma(w_t, v_t)$. Together, these two negative terms contract
509 the Lyapunov potential $\Psi_{\alpha}(w_t, v_t) = (\Phi_D(w_t) - \Phi_D^*) + \alpha \Gamma(w_t, v_t)$, yielding a multiplicative

510 shrinkage factor $(1 - c\gamma)\Psi_\alpha(w_t, v_t)$ with $c := \min\{\mu_{\text{PL}}/2, \rho/2\}$. Furthermore, the positive
 511 $\gamma\Gamma$ pieces produced by the algebra consist of

$$\alpha \cdot \frac{\rho}{2} \gamma \Gamma(w_t, v_t) \quad \text{and} \quad \frac{L^2}{\rho} \gamma \Gamma(w_t, v_t).$$

512 Together with the dual drift they appear as

$$-2\alpha\rho\gamma\Gamma + \alpha \cdot \frac{\rho}{2} \gamma \Gamma + \frac{L^2}{\rho} \gamma \Gamma = -\frac{3}{2}\alpha\rho\gamma\Gamma + \frac{L^2}{\rho} \gamma \Gamma.$$

513 Choose α so that $\frac{L^2}{\rho} \leq \frac{\alpha\rho}{2}$ (i.e. $\alpha \geq 2L^2/\rho^2$). Then

$$-\frac{3}{2}\alpha\rho\gamma\Gamma + \frac{L^2}{\rho} \gamma \Gamma \leq -\frac{\alpha\rho}{2} \gamma \Gamma,$$

514 so these $\gamma\Gamma$ remainders are dominated by the dual drift and absorbed into the contraction factor.
 515 A simple sufficient standing choice is $\alpha \geq 4L^2/\rho^2$.

516 (ii) The mismatch bounds, ghost-correction terms, and gradient–difference terms contribute constants
 517 times γG^2 . All such pieces are nonnegative and can be grouped into a single term
 518 $C_{\text{lin}} G^2 \gamma$.

519 (iii) Smoothness corrections (e.g. $(L_\Phi + L)/2 \gamma^2 G^2, 2L \gamma^2 G^2, L \gamma^2 G^2$) contribute constants times
 520 $\gamma^2 G^2$. These are also nonnegative and can be collected into $A \gamma^2$.

521 Putting the three groups together, the one-step recursion takes the form

$$\mathbb{E}[\Psi_\alpha(w^+, v^+) | \mathcal{F}_t] \leq (1 - c\gamma) \Psi_\alpha(w_t, v_t) + C_{\text{lin}} G^2 \gamma + A \gamma^2,$$

522 which is exactly the recursion **(P)** used in the stability analysis.

523 **(c) Two-run recursion.** Recall $\Xi_t := \Psi_{\alpha, \mathcal{D}}(w_t, v_t) + \Psi_{\alpha, \mathcal{D}^{(i)}}(w'_t, v'_t)$ and $D_t := \|w_t - w'_t\| +$
 524 $\|v_t - v'_t\|$. Apply the one-run bound from part (b) to \mathcal{D} and to $\mathcal{D}^{(i)}$ (same stepsize; shared index i_t),
 525 then sum:

$$\begin{aligned} \mathbb{E}[\Xi_{t+1} | \mathcal{F}_t] &\leq (1 - c\eta_t) \Xi_t + A \eta_t^2 + C_{\text{lin}} G^2 \eta_t \\ &\quad + \underbrace{\eta_t \left(\langle \nabla \Phi_{\mathcal{D}}(w_t) - \nabla \Phi_{\mathcal{D}^{(i)}}(w'_t), g_t^w \rangle - \langle \nabla_w F_{\mathcal{D}}(w_t, v_t) - \nabla_w F_{\mathcal{D}^{(i)}}(w'_t, v'_t), g_t^w \rangle \right)}_{=: \text{Leak}_t^{(w)}} \\ &\quad + \eta_t \left(\mathbf{1}\{i_t = i\} + \mathbf{1}\{\hat{i}_t = i\} \right) \tilde{B}_0 G^2. \end{aligned} \tag{13}$$

526 Bounding the leak term by Cauchy–Schwarz, joint L -smoothness, and Lemma 10 gives

$$\begin{aligned} |\text{Leak}_t^{(w)}| &= \eta_t \left| \langle \nabla \Phi_{\mathcal{D}}(w_t) - \nabla \Phi_{\mathcal{D}^{(i)}}(w'_t), g_t^w \rangle - \langle \nabla_w F_{\mathcal{D}}(w_t, v_t) - \nabla_w F_{\mathcal{D}^{(i)}}(w'_t, v'_t), g_t^w \rangle \right| \\ &\leq \eta_t \left(\left| \langle \nabla \Phi_{\mathcal{D}}(w_t) - \nabla \Phi_{\mathcal{D}^{(i)}}(w'_t), g_t^w \rangle \right| + \left| \langle \nabla_w F_{\mathcal{D}}(w_t, v_t) - \nabla_w F_{\mathcal{D}^{(i)}}(w'_t, v'_t), g_t^w \rangle \right| \right) \\ &\leq \eta_t \left(\left\| \nabla \Phi_{\mathcal{D}}(w_t) - \nabla \Phi_{\mathcal{D}^{(i)}}(w'_t) \right\| \|g_t^w\| + \left\| \nabla_w F_{\mathcal{D}}(w_t, v_t) - \nabla_w F_{\mathcal{D}^{(i)}}(w'_t, v'_t) \right\| \|g_t^w\| \right) \\ &\hspace{15em} \text{(Cauchy-Schwarz)} \\ &\leq \eta_t \left(L \|w_t - w'_t\| \|g_t^w\| + L_\Phi \|w_t - w'_t\| \|g_t^w\| + 2L \|v_t - v'_t\| \|g_t^w\| \right) \\ &\hspace{15em} \text{(Joint } L\text{-smoothness of gradients)} \\ &\leq \eta_t G \left((2L + 2L_\Phi) \|w_t - w'_t\| + 2L \|v_t - v'_t\| + \frac{2G}{n} \left(2 + \frac{L}{\rho} \right) \right). \quad \text{(Lemma 10)} \end{aligned}$$

527 hence, with $D_t = \|w_t - w'_t\| + \|v_t - v'_t\|$,

$$\text{Leak}_t := \text{Leak}_t^{(w)} \leq C_{\text{leak}} G \eta_t D_t + C_* \frac{G^2}{n} \eta_t, \quad C_{\text{leak}} := 2(2L + L_\Phi), \quad C_* := 4 \left(2 + \frac{L}{\rho} \right).$$

528 Therefore

$$\mathbb{E}[\Xi_{t+1} \mid \mathcal{F}_t] \leq (1 - c\eta_t) \Xi_t + A\eta_t^2 + C_{\text{leak}} G \eta_t D_t + C_* \frac{G^2}{n} \eta_t + C_{\text{lin}} G^2 \eta_t + \eta_t (\mathbf{1}\{i_t = i\} + \mathbf{1}\{\hat{i}_t = i\}) \tilde{B}_0 G^2, \quad (14)$$

529 which matches (6) up to explicit constants.

530 *Weak distance recursion.* A single SGDA step gives

$$\|w_{t+1} - w'_{t+1}\| \leq \|w_t - w'_t\| + \eta_t \|\nabla_w f(w_t, v_t; z_{i_t}) - \nabla_w f(w'_t, v'_t; z_{i_t})\| + \eta_t \mathbf{1}\{i_t = i\} \cdot 2G,$$

531 and similarly for v (with ascent). Joint L -smoothness yields

$$\|\nabla_w f(w_t, v_t; z) - \nabla_w f(w'_t, v'_t; z)\| \leq L(\|w_t - w'_t\| + \|v_t - v'_t\|), \quad \|\nabla_v f(\cdot) - \nabla_v f(\cdot)\| \leq L(\|w_t - w'_t\| + \|v_t - v'_t\|).$$

532 Hence

$$\mathbb{E}[D_{t+1} \mid \mathcal{F}_t] \leq (1 + \kappa\eta_t) D_t + \eta_t \mathbf{1}\{i_t = i\} \tilde{B}_1 G + \eta_t \Upsilon_t,$$

533 with $\kappa := 2L$ and $\tilde{B}_1 := 8$.

534 To bound the term Υ_t , which represents the expected difference in stochastic gradients, we first
535 analyze the difference of the full-batch gradients. We decompose this difference for the primal and
536 dual variables separately.

537 For the primal variable, we use the triangle inequality to introduce the value function Φ and the
538 mismatch term $\Delta_t := \nabla_w F_{\mathcal{D}}(w_t, v_t) - \nabla \Phi_{\mathcal{D}}(w_t)$:

$$\begin{aligned} \|\nabla_w F_{\mathcal{D}}(w_t, v_t) - \nabla_w F_{\mathcal{D}^{(i)}}(w'_t, v'_t)\| &\leq \underbrace{\|\nabla \Phi_{\mathcal{D}}(w_t) - \nabla \Phi_{\mathcal{D}^{(i)}}(w'_t)\|}_{\leq L_{\Phi} \|w_t - w'_t\| + \frac{2G}{n} (1 + \frac{L}{\rho}) \text{ by Lem. 10}} \\ &+ \underbrace{\|\Delta_t\|}_{\leq L\sqrt{2\Gamma_t/\rho} \text{ by Lem. 8}} + \underbrace{\|\Delta'_t\|}_{\leq L\sqrt{2\Gamma'_t/\rho} \text{ by Lem. 8}} \end{aligned}$$

539 For the dual variable, we use joint L -smoothness and the single-sample replacement identity:

$$\begin{aligned} \|\nabla_v F_{\mathcal{D}}(w_t, v_t) - \nabla_v F_{\mathcal{D}^{(i)}}(w'_t, v'_t)\| &\leq \|\nabla_v F_{\mathcal{D}}(w_t, v_t) - \nabla_v F_{\mathcal{D}}(w'_t, v'_t)\| + \|\nabla_v F_{\mathcal{D}}(w'_t, v'_t) - \nabla_v F_{\mathcal{D}^{(i)}}(w'_t, v'_t)\| \\ &\leq L(\|w_t - w'_t\| + \|v_t - v'_t\|) + \frac{2G}{n} = LD_t + \frac{2G}{n}. \end{aligned}$$

540 Combining the bounds on the primal and dual components gives a complete bound on the full-batch
541 gradient difference. Summing the two inequalities and using $D_t \geq \|w_t - w'_t\|$, we have:

$$\|\nabla F_{\mathcal{D}}(w_t, v_t) - \nabla F_{\mathcal{D}^{(i)}}(w'_t, v'_t)\|_1 \leq (L_{\Phi} + L)D_t + L\sqrt{2/\rho}(\sqrt{\Gamma_t} + \sqrt{\Gamma'_t}) + \frac{2G}{n} \left(2 + \frac{L}{\rho}\right).$$

542 We then convert the geometric quantities on the right-hand side (D_t and Γ_t) into the Lyapunov
543 potential Ξ_t . First, using the triangle inequality along with Lemmas 11 and 12, we bound the distance
544 D_t :

$$\begin{aligned} D_t &= \|w_t - w'_t\| + \|v_t - v'_t\| \\ &\leq (\|w_t - x_D^*\| + \|v_t - v_D^*\|) + (\|x_D^* - x_{D^{(i)}}^*\| + \|v_D^* - v_{D^{(i)}}^*\|) + (\|x_{D^{(i)}}^* - w'_t\| + \|v_{D^{(i)}}^* - v'_t\|) \\ &\leq C_{\text{dist}} \sqrt{\Psi_{\alpha, \mathcal{D}}(w_t, v_t)} + \frac{S_{\text{sens}}}{n} + C_{\text{dist}} \sqrt{\Psi_{\alpha, \mathcal{D}^{(i)}}(w'_t, v'_t)} \\ &\leq C_{\text{dist}} \left(\sqrt{\Psi_{\alpha, \mathcal{D}}} + \sqrt{\Psi_{\alpha, \mathcal{D}^{(i)}}} \right) + \frac{S_{\text{sens}}}{n} \leq \sqrt{2} C_{\text{dist}} \sqrt{\Xi_t} + \frac{S_{\text{sens}}}{n}. \end{aligned}$$

545 Second, from the definition of the potential, we have $\sqrt{\Gamma_t} + \sqrt{\Gamma'_t} \leq 2\sqrt{\Xi_t/\alpha}$. Substituting these
546 into our main inequality gives:

$$\|\nabla F_{\mathcal{D}}(w_t, \cdot) - \nabla F_{\mathcal{D}^{(i)}}(w'_t, \cdot)\|_1 \leq (L_{\Phi} + L) \left(\sqrt{2} C_{\text{dist}} \sqrt{\Xi_t} + \frac{S_{\text{sens}}}{n} \right) + L\sqrt{2/\rho} \left(2\sqrt{\Xi_t/\alpha} \right) + O(G/n).$$

547 The terms involving $\sqrt{\Xi_t}$ can be collected and absorbed into a single term $a\sqrt{\Xi_t}$, where a is a constant
548 that depends on the problem parameters (e.g., $a := 2C_{\text{dist}}(L + L_{\Phi})$ serves as a valid, convenient

549 upper bound). The remaining terms are of order $O(1/n)$. Therefore, we state the final bound on Υ_t
 550 as:

$$\Upsilon_t \leq 2C_{\text{dist}}(L + L_{\Phi}) \sqrt{\Xi_t} + \frac{2G}{n} \left(1 + \frac{L}{\rho}\right),$$

551 and therefore

$$\mathbb{E}[D_{t+1} \mid \mathcal{F}_t] \leq (1 + \kappa\eta_t) D_t + a\eta_t \sqrt{\Xi_t} + \eta_t \mathbf{1}\{i_t = i\} \tilde{B}_1 G + \eta_t \frac{2G}{n} \left(1 + \frac{L}{\rho}\right), \quad a := 2C_{\text{dist}}(L + L_{\Phi}). \quad (15)$$

552

□

553 **Lemma 14** (One-step ascent for L -smooth, ρ -strongly concave g). *Let $g : \mathbb{R}^d \rightarrow \mathbb{R}$ be L -smooth
 554 and ρ -strongly concave, and let $v^+ = v + \gamma \nabla g(v)$ with $\gamma \geq 0$. Denote $\theta(v) := g(v^*) - g(v)$, where
 555 $v^* \in \arg \max g$. Then*

$$\theta(v^+) \leq (1 - 2\rho\gamma + \rho L\gamma^2) \theta(v).$$

556 *Consequently,*

$$g(v^*) - g(v^+) \leq (1 - 2\rho\gamma + \rho L\gamma^2) [g(v^*) - g(v)].$$

557 *Proof.* For any L -smooth differentiable function h one has

$$h(y) \geq h(x) + \langle \nabla h(x), y - x \rangle - \frac{L}{2} \|y - x\|^2.$$

558 Apply this to $h = g$, $x = v$, $y = v^+ = v + \gamma \nabla g(v)$:

$$g(v^+) \geq g(v) + \gamma \|\nabla g(v)\|^2 - \frac{L}{2} \gamma^2 \|\nabla g(v)\|^2.$$

559 Rearranging gives

$$g(v^*) - g(v^+) \leq g(v^*) - g(v) - \left(\gamma - \frac{L}{2}\gamma^2\right) \|\nabla g(v)\|^2. \quad (16)$$

560 Since g is ρ -strongly concave, $f := -g$ is ρ -strongly convex; hence f satisfies the Polyak–Łojasiewicz
 561 (PL) inequality

$$\frac{1}{2} \|\nabla f(x)\|^2 \geq \rho(f(x) - f^*),$$

562 with the *same* constant ρ [Karimi et al., 2016]. Translating back to g gives

$$\frac{1}{2} \|\nabla g(v)\|^2 \geq \rho(g(v^*) - g(v)) = \rho\theta(v). \quad (17)$$

563 *Step 3 (Combine).* Insert (17) into (16):

$$\theta(v^+) \leq \theta(v) - 2\rho\gamma \left(1 - \frac{L}{2}\gamma\right) \theta(v) = (1 - 2\rho\gamma + \rho L\gamma^2) \theta(v),$$

564

□

565 **Lemma 15** (Damping the weak distance recursion). *Let $S_t := \sum_{s=0}^{t-1} \eta_s$ and define the damped
 566 distance $\tilde{D}_t := e^{-2LS_t} D_t$ from (7). Set $\lambda := 2L$ for this lemma. Then, taking expectations and
 567 averaging over i ,*

$$\tilde{D}_{t+1} \leq \tilde{D}_t + a\eta_t e^{-2LS_{t+1}} \mathbb{E}[\sqrt{\Xi_t}] + \frac{1}{n} \left(2\tilde{B}_1 G + 2G \left(1 + \frac{L}{\rho}\right)\right) \eta_t e^{-\lambda S_{t+1}}.$$

568 *Consequently, summing from $t = 0$ to $T - 1$,*

$$\bar{D}_T \leq a \sum_{t=0}^{T-1} \eta_t e^{-2L \sum_{s=t+1}^{T-1} \eta_s} \mathbb{E}[\sqrt{\Xi_t}] + \frac{1}{\lambda n} \left(2\tilde{B}_1 G + 2G \left(1 + \frac{L}{\rho}\right)\right). \quad (18)$$

569 *Proof.* Start from the weak distance recursion (7):

$$\mathbb{E}[D_{t+1} \mid \mathcal{F}_t] \leq (1 + 2L\eta_t) D_t + a\eta_t \sqrt{\Xi_t} + \eta_t \mathbf{1}\{i_t = i\} \tilde{B}_1 G + \eta_t \frac{2G}{n} \left(1 + \frac{L}{\rho}\right).$$

570 Multiply both sides by $e^{-2LS_{t+1}} = e^{-2L(S_t + \eta_t)}$ to get

$$e^{-2LS_{t+1}} \mathbb{E}[D_{t+1} | \mathcal{F}_t] \leq e^{-2LS_t} (1 + 2L\eta_t) e^{-2L\eta_t} D_t + a\eta_t e^{-2LS_{t+1}} \sqrt{\Xi_t} \\ + \eta_t e^{-2LS_{t+1}} \mathbf{1}\{i_t = i\} \tilde{B}_1 G + \eta_t e^{-2LS_{t+1}} \frac{2G}{n} \left(1 + \frac{L}{\rho}\right).$$

571 Since $(1+x)e^{-x} \leq 1$ for all $x \geq 0$, the first term is $\leq e^{-2LS_t} D_t = \tilde{D}_t$. Taking total expectation
572 and then averaging over the replacement index i , we make explicit that the hit $\mathbf{1}\{i_t = i\}$ contributes
573 to *both* coordinates in $D_t = \|w_t - w'_t\| + \|v_t - v'_t\|$. For a fixed i ,

$$\mathbb{E}[D_{t+1} | \mathcal{F}_t] \leq (1 + 2L\eta_t) D_t + a\eta_t \sqrt{\Xi_t} + \eta_t \mathbf{1}\{i_t = i\} c_w G + \eta_t \mathbf{1}\{i_t = i\} c_v G + \eta_t \frac{2G}{n} \left(1 + \frac{L}{\rho}\right) \\ \leq (1 + 2L\eta_t) D_t + a\eta_t \sqrt{\Xi_t} + \eta_t \mathbf{1}\{i_t = i\} (2\tilde{B}_1) G + \eta_t \frac{2G}{n} \left(1 + \frac{L}{\rho}\right),$$

574 where we bundle constants so that $c_w = c_v = \tilde{B}_1$. Multiplying both sides by $e^{-2LS_{t+1}}$ with
575 $S_t := \sum_{s=0}^{t-1} \eta_s$ and using $(1 + 2L\eta_t) e^{-2L\eta_t} \leq 1$ yields

$$\mathbb{E}[\tilde{D}_{t+1} | \mathcal{F}_t] \leq \tilde{D}_t + a\eta_t e^{-2LS_{t+1}} \sqrt{\Xi_t} + \eta_t e^{-2LS_{t+1}} \mathbf{1}\{i_t = i\} (2\tilde{B}_1) G + \eta_t e^{-2LS_{t+1}} \frac{2G}{n} \left(1 + \frac{L}{\rho}\right),$$

576 where $\tilde{D}_t := e^{-2LS_t} D_t$. Taking total expectation and then averaging over i (so that $\mathbb{E}[\mathbf{1}\{i_t = i\}] =$
577 $1/n$) gives

$$\tilde{D}_{t+1} \leq \tilde{D}_t + a\eta_t e^{-2LS_{t+1}} \mathbb{E}[\sqrt{\Xi_t}] + \frac{1}{n} \left(2\tilde{B}_1 G + 2G \left(1 + \frac{L}{\rho}\right)\right) \eta_t e^{-2LS_{t+1}}.$$

578 Summing the inequality from $t = 0$ to $T-1$ and noting that both runs start from the same initialization
579 $(w_0, v_0) = (w'_0, v'_0)$, we have $D_0 = 0$ and hence $\tilde{D}_0 = 0$ and $\bar{D}_0 = 0$. Therefore,

$$\tilde{D}_T - \tilde{D}_0 = \tilde{D}_T \leq a \sum_{t=0}^{T-1} \eta_t e^{-2LS_{t+1}} \mathbb{E}[\sqrt{\Xi_t}] + \frac{1}{n} \left(2\tilde{B}_1 G + 2G \left(1 + \frac{L}{\rho}\right)\right) \sum_{t=0}^{T-1} \eta_t e^{-2LS_{t+1}}.$$

580 Applying Lemma 16 with $\gamma = \lambda = 2L$ to the last sum and noting that $e^{2LS_T} e^{-2LS_{t+1}} =$
581 $e^{2L \sum_{s=t+1}^{T-1} \eta_s}$ yields

$$\bar{D}_T = e^{2LS_T} \tilde{D}_T \leq a \sum_{t=0}^{T-1} \eta_t e^{-2L \sum_{s=t+1}^{T-1} \eta_s} \mathbb{E}[\sqrt{\Xi_t}] + \frac{1}{\lambda n} \left(2\tilde{B}_1 G + 2G \left(1 + \frac{L}{\rho}\right)\right),$$

582 which is (18). □

583 **Lemma 16.** Let (η_t) satisfy the Robbins–Monro conditions. For $S_t := \sum_{s=0}^{t-1} \eta_s$ and any $\gamma > 0$, any
584 $T \geq 1$, set $\eta_{\max, T} := \max_{0 \leq t < T} \eta_t$. Then

$$\sum_{t=0}^{T-1} \eta_t e^{-\gamma \sum_{s=t+1}^{T-1} \eta_s} \leq \frac{e^{\gamma \eta_{\max, T}}}{\gamma} (1 - e^{-\gamma S_T}) \leq \frac{e^{\gamma \eta_{\max, T}}}{\gamma},$$

585 and

$$\sum_{t=0}^{T-1} \eta_t^2 e^{-\gamma \sum_{s=t+1}^{T-1} \eta_s} \xrightarrow{T \rightarrow \infty} 0.$$

586 In particular, under Condition (S1) and with $\gamma = 2L$, one has $\eta_{\max, T} \leq 1/(4L)$ and therefore

$$\sum_{t=0}^{T-1} \eta_t e^{-2L \sum_{s=t+1}^{T-1} \eta_s} \leq \frac{e^{1/2}}{2L}.$$

587 *Proof.* Let $S_t := \sum_{s=0}^{t-1} \eta_s$ and note $S_t \uparrow \infty$ under Robbins–Monro. For the first display, fix t and
588 any $u \in [S_t, S_{t+1}]$. Then

$$e^{-\gamma(S_T - S_{t+1})} = e^{-\gamma(S_T - u)} e^{\gamma(S_{t+1} - u)} \leq e^{\gamma \eta_{\max, T}} e^{-\gamma(S_T - u)},$$

589 since $0 \leq S_{t+1} - u \leq S_{t+1} - S_t = \eta_t \leq \eta_{\max, T}$. Hence

$$\eta_t e^{-\gamma(S_T - S_{t+1})} \leq e^{\gamma\eta_{\max, T}} \int_{S_t}^{S_{t+1}} e^{-\gamma(S_T - u)} du.$$

590 Summing over $t = 0, \dots, T-1$ and using $\sum_t \int_{S_t}^{S_{t+1}} (\cdot) = \int_0^{S_T} (\cdot)$ gives

$$\sum_{t=0}^{T-1} \eta_t e^{-\gamma \sum_{s=t+1}^{T-1} \eta_s} \leq e^{\gamma\eta_{\max, T}} \int_0^{S_T} e^{-\gamma(S_T - u)} du = \frac{e^{\gamma\eta_{\max, T}}}{\gamma} (1 - e^{-\gamma S_T}) \leq \frac{e^{\gamma\eta_{\max, T}}}{\gamma}.$$

591 For the second display, fix $\varepsilon > 0$. Since $\sum_t \eta_t^2 < \infty$, pick T_0 with $\sum_{t \geq T_0} \eta_t^2 < \varepsilon$. Then split

$$\sum_{t=0}^{T-1} \eta_t^2 e^{-\gamma \sum_{s=t+1}^{T-1} \eta_s} = \sum_{t=0}^{T_0-1} \eta_t^2 e^{-\gamma(S_T - S_{t+1})} + \sum_{t=T_0}^{T-1} \eta_t^2 e^{-\gamma(S_T - S_{t+1})} \leq e^{-\gamma(S_T - S_{T_0})} \sum_{t=0}^{T_0-1} \eta_t^2 + \varepsilon.$$

592 Let $T \rightarrow \infty$ to send the first term to 0 (since $S_T \rightarrow \infty$), and then let $\varepsilon \downarrow 0$. Finally, under
 593 Condition (S1) and $\gamma = 2L$, we have $\eta_{\max, T} \leq 1/(4L)$, hence $e^{\gamma\eta_{\max, T}} \leq e^{1/2}$, which yields the
 594 stated corollary. \square

595 **Lemma 17** (Closing the potential recursion). *After averaging (6) over i ,*

$$\bar{\Xi}_{t+1} \leq \left(1 - \frac{c}{2}\eta_t\right)\bar{\Xi}_t + A\eta_t^2 + \frac{b^2}{2c}\eta_t + \frac{\tilde{H}}{n}\eta_t, \quad (19)$$

596 with $b := 2C_{\text{leak}}C_{\text{dist}}$ and $\tilde{H} := 2\tilde{B}_0G^2 + C_{\text{leak}}S_{\text{sens}} + C_*G$, where

$$S_{\text{sens}} := 2G \left(\frac{(1 + L/\rho)^2}{\sqrt{\mu_{\text{PL}}\mu_{\text{QG}}}} + \frac{1}{\rho} \right)$$

597 is the constant from Lemma 12.

598 *Proof.* Starting from (6),

$$\mathbb{E}[\bar{\Xi}_{t+1} \mid \mathcal{F}_t] \leq (1 - c\eta_t)\bar{\Xi}_t + A\eta_t^2 + C_{\text{leak}}\eta_t D_t + C_*\frac{G}{n}\eta_t + \eta_t(\mathbf{1}\{i_t = i\} + \mathbf{1}\{\hat{i}_t = i\})\tilde{B}_0G^2.$$

599 By Lemma 11, for each run $\|w_t - x_D^*\| + \|v_t - v_D^*\| \leq C_{\text{dist}}\sqrt{\Psi_{\alpha, \mathcal{D}}(w_t, v_t)}$ and analogously for the
 600 primed run. Using the triangle inequality and Lemma 12 for the cross-dataset saddle shift, we have

$$\begin{aligned} D_t &= \|w_t - w'_t\| + \|v_t - v'_t\| \\ &\leq \|w_t - x_D^*\| + \|x_D^* - x_{D^{(i)}}^*\| + \|x_{D^{(i)}}^* - w'_t\| + \|v_t - v_D^*\| + \|v_D^* - v_{D^{(i)}}^*\| - \|v_{D^{(i)}}^* - v'_t\| \\ &\quad \text{(Triangle inequality)} \\ &\leq C_{\text{dist}} \left(\sqrt{\Psi_{\alpha, \mathcal{D}}(w_t, v_t)} + \sqrt{\Psi_{\alpha, \mathcal{D}^{(i)}}(w'_t, v'_t)} \right) + \|x_D^* - x_{D^{(i)}}^*\| - \|v_{D^{(i)}}^* - v'_t\| \\ &\quad \text{(Lemma 12)} \\ &\leq C_{\text{dist}} \left(\sqrt{\Psi_{\alpha, \mathcal{D}}(w_t, v_t)} + \sqrt{\Psi_{\alpha, \mathcal{D}^{(i)}}(w'_t, v'_t)} \right) + \frac{S_{\text{sens}}}{n} \\ &\quad \text{(Lemma 11)} \end{aligned}$$

601 Since $\sqrt{\Psi_{\alpha, \mathcal{D}}} \leq \sqrt{\bar{\Xi}_t}$ and $\sqrt{\Psi_{\alpha, \mathcal{D}^{(i)}}} \leq \sqrt{\bar{\Xi}_t}$,

$$D_t \leq 2C_{\text{dist}}\sqrt{\bar{\Xi}_t} + \frac{S_{\text{sens}}}{n}.$$

602 Plug this into (6):

$$\mathbb{E}[\bar{\Xi}_{t+1} \mid \mathcal{F}_t] \leq (1 - c\eta_t)\bar{\Xi}_t + A\eta_t^2 + \underbrace{2C_{\text{leak}}C_{\text{dist}}}_{=: b}\eta_t\sqrt{\bar{\Xi}_t} + \frac{C_{\text{leak}}S_{\text{sens}}}{n}\eta_t + C_*\frac{G}{n}\eta_t + \eta_t(\mathbf{1}\{i_t = i\} + \mathbf{1}\{\hat{i}_t = i\})\tilde{B}_0G^2.$$

603 Apply the inequality $uv \leq \frac{\gamma}{2}u^2 + \frac{1}{2\gamma}v^2$ with $u = \sqrt{\eta_t\bar{\Xi}_t}$, $v = b\sqrt{\eta_t}$, and $\gamma = c$ to the mixed term:

$$b\eta_t\sqrt{\bar{\Xi}_t} \leq \frac{c}{2}\eta_t\bar{\Xi}_t + \frac{b^2}{2c}\eta_t.$$

604 Therefore,

$$\mathbb{E}[\Xi_{t+1} \mid \mathcal{F}_t] \leq \left(1 - \frac{c}{2}\eta_t\right)\Xi_t + A\eta_t^2 + \frac{b^2}{2c}\eta_t + \frac{C_{\text{leak}}S_{\text{sens}}}{n}\eta_t + C_*\frac{G}{n}\eta_t + \eta_t(\mathbf{1}\{i_t = i\} + \mathbf{1}\{\hat{i}_t = i\})\tilde{B}_0G^2.$$

605 Taking total expectation and averaging over i (both indicators have mean $1/n$) gives

$$\bar{\Xi}_{t+1} \leq \left(1 - \frac{c}{2}\eta_t\right)\bar{\Xi}_t + A\eta_t^2 + \frac{b^2}{2c}\eta_t + \frac{1}{n}\left(2\tilde{B}_0G^2 + C_{\text{leak}}S_{\text{sens}} + C_*G\right)\eta_t.$$

606 With $\tilde{H} := 2\tilde{B}_0G^2 + C_{\text{leak}}S_{\text{sens}} + C_*G$ and $b = 2C_{\text{leak}}C_{\text{dist}}$ this is (19). \square

607 **Lemma 18** (Bounded weighted sum of potentials). *For any $\theta \in (\frac{1}{2}, 1)$ and $S_t := \sum_{s=0}^{t-1} \eta_s$,*

$$\begin{aligned} \sum_{t=0}^{T-1} e^{-\theta c \sum_{s=t+1}^{T-1} \eta_s} \bar{\Xi}_t &\leq \frac{2}{1 - e^{-(\theta - \frac{1}{2})c \eta_{\min, T}}} e^{-\frac{c}{2} \sum_{s=0}^{T-1} \eta_s} \bar{\Xi}_0 \\ &+ \frac{2}{1 - e^{-(\theta - \frac{1}{2})c \eta_{\min, T}}} A \sum_{t=0}^{T-1} \eta_t^2 e^{-\frac{c}{2} \sum_{s=t+1}^{T-1} \eta_s} + \frac{8}{c(1 - e^{-(\theta - \frac{1}{2})c \eta_{\min, T}})} \left(\frac{b^2}{2c} + \frac{\tilde{H}}{n}\right). \end{aligned}$$

608 *In particular, since $c \eta_{\min, T} \leq \frac{1}{2}$ (by Condition (S1)) and hence $1 - e^{-x} \geq x/2$ for $x \in [0, \frac{1}{2}]$, we*
609 *have*

$$\begin{aligned} \sum_{t=0}^{T-1} e^{-\theta c \sum_{s=t+1}^{T-1} \eta_s} \bar{\Xi}_t &\leq \frac{4}{(\theta - \frac{1}{2})c \eta_{\min, T}} e^{-\frac{c}{2} \sum_{s=0}^{T-1} \eta_s} \bar{\Xi}_0 \\ &+ \frac{4}{(\theta - \frac{1}{2})c \eta_{\min, T}} A \sum_{t=0}^{T-1} \eta_t^2 e^{-\frac{c}{2} \sum_{s=t+1}^{T-1} \eta_s} + \frac{16}{(\theta - \frac{1}{2})c^2 \eta_{\min, T}} \left(\frac{b^2}{2c} + \frac{\tilde{H}}{n}\right). \end{aligned}$$

610 *Proof.* Start from (19) and unroll using $1 - \frac{c}{2}\eta_k \leq e^{-\frac{c}{2}\eta_k}$:

$$\bar{\Xi}_t \leq e^{-\frac{c}{2}S_t} \bar{\Xi}_0 + \sum_{k=0}^{t-1} e^{-\frac{c}{2}(S_t - S_{k+1})} (A\eta_k^2 + q\eta_k), \quad q := \frac{b^2}{2c} + \frac{\tilde{H}}{n}.$$

611 Multiply by the weight $w_t := e^{-\theta c(S_T - S_{t+1})}$ and sum over $t = 0, \dots, T-1$:

$$\sum_{t=0}^{T-1} w_t \bar{\Xi}_t \leq \underbrace{\sum_{t=0}^{T-1} w_t e^{-\frac{c}{2}S_t} \bar{\Xi}_0}_I + A \underbrace{\sum_{t=0}^{T-1} \sum_{k=0}^{t-1} w_t e^{-\frac{c}{2}(S_t - S_{k+1})} \eta_k^2}_{II} + q \underbrace{\sum_{t=0}^{T-1} \sum_{k=0}^{t-1} w_t e^{-\frac{c}{2}(S_t - S_{k+1})} \eta_k}_{III}.$$

612 Throughout, Condition (S1) implies $c\eta_t \leq \frac{1}{2}$ (since $c \leq \rho/2$ and $\eta_t \leq 1/\rho$), hence $1 - e^{-x} \geq x/2$
613 for $x \in [0, \frac{1}{2}]$.

614 *Claim A.* For $\theta \in (\frac{1}{2}, 1)$,

$$I = \sum_{t=0}^{T-1} e^{-\theta c(S_T - S_{t+1})} e^{-\frac{c}{2}S_t} \leq \frac{2e^{-\frac{c}{2}S_T}}{1 - \exp(-(\theta - \frac{1}{2})c \eta_{\min, T})} \leq \frac{4}{(\theta - \frac{1}{2})c \eta_{\min, T}} e^{-\frac{c}{2}S_T},$$

615 where $\eta_{\min, T} := \min_{0 \leq t < T} \eta_t$.

616 *Proof of Claim A.* Using $S_t = S_{t+1} - \eta_t$ and $S_T = S_{t+1} + (S_T - S_{t+1})$,

$$e^{-\theta c(S_T - S_{t+1})} e^{-\frac{c}{2}S_t} = e^{-\frac{c}{2}S_T} e^{-(\theta - \frac{1}{2})c(S_T - S_{t+1})} e^{\frac{c}{2}\eta_t}.$$

617 Since $c\eta_t \leq \frac{1}{2}$, $e^{\frac{c}{2}\eta_t} \leq e^{1/4} \leq 2$. Therefore,

$$I \leq 2e^{-\frac{c}{2}S_T} \sum_{t=0}^{T-1} e^{-(\theta - \frac{1}{2})c(S_T - S_{t+1})}.$$

618 Because $S_T - S_{t+1} \geq (T-1-t)\eta_{\min,T}$, the sum is bounded by a geometric series, giving the first
619 display; the second follows from $1 - e^{-x} \geq x/2$ for $x \in [0, 1]$. \triangle

620 *Claim B.* Let $\alpha := (\theta - \frac{1}{2})c > 0$. For any fixed $k \in \{0, \dots, T-2\}$,

$$\sum_{t=k+1}^{T-1} w_t e^{-\frac{c}{2}(S_t - S_{k+1})} \leq \frac{2}{1 - e^{-\alpha\eta_{\min,T}}} e^{-\frac{c}{2}(S_T - S_{k+1})} \leq \frac{4}{(\theta - \frac{1}{2})c\eta_{\min,T}} e^{-\frac{c}{2}(S_T - S_{k+1})}.$$

621 *Proof of Claim B.* As above,

$$w_t e^{-\frac{c}{2}(S_t - S_{k+1})} = e^{-\frac{c}{2}(S_T - S_{k+1})} e^{-\alpha(S_T - S_{t+1})} e^{\frac{c}{2}\eta_t} \leq 2 e^{-\frac{c}{2}(S_T - S_{k+1})} e^{-\alpha(S_T - S_{t+1})}.$$

622 Let $I_t := \int_{S_{t+1}}^{S_{t+2}} e^{-\alpha(S_T - u)} du$. A direct computation gives

$$I_t = \frac{1}{\alpha} e^{-\alpha(S_T - S_{t+1})} (1 - e^{-\alpha\eta_{t+1}}) \Rightarrow e^{-\alpha(S_T - S_{t+1})} = \frac{\alpha}{1 - e^{-\alpha\eta_{t+1}}} I_t \leq \frac{\alpha}{1 - e^{-\alpha\eta_{\min,T}}} I_t,$$

623 since $\eta_{t+1} \geq \eta_{\min,T}$. Summing over $t = k+1, \dots, T-1$ and using $\sum_{t=k+1}^{T-1} I_t =$
624 $\int_{S_{k+1}}^{S_T} e^{-\alpha(S_T - u)} du \leq 1/\alpha$ yields the claim; the explicit bound uses $1 - e^{-x} \geq x/2$ with
625 $x = \alpha\eta_{\min,T} \leq \frac{1}{2}$. \triangle

626 With Claim A and $1 - e^{-x} > x/2$,

$$I \leq \frac{2}{1 - e^{-(\theta - \frac{1}{2})c\eta_{\min,T}}} e^{-\frac{c}{2}S_T},$$

627 and, applying Claim B for each fixed k inside the double sum and then summing in t ,

$$\begin{aligned} \text{II} &= \sum_{t=0}^{T-1} \sum_{k=0}^{t-1} w_t e^{-\frac{c}{2}(S_t - S_{k+1})} \eta_k^2 \\ &= \sum_{k=0}^{T-2} \sum_{t=k+1}^{T-1} w_t e^{-\frac{c}{2}(S_t - S_{k+1})} \eta_k^2 && \text{(Fubini)} \\ &\leq \sum_{k=0}^{T-2} \frac{2}{1 - e^{-\alpha\eta_{\min,T}}} e^{-\frac{c}{2}(S_T - S_{k+1})} \eta_k^2 && \text{(Claim B)} \end{aligned}$$

628 For III, combine Claim B with the kernel bound from Lemma 16

$$\sum_{k=0}^{T-2} \eta_k e^{-\frac{c}{2}(S_T - S_{k+1})} \leq \frac{2 e^{\frac{c}{2}\eta_{\max,T}}}{c} \leq \frac{4}{c} \quad (\text{since } c\eta_{\max,T} \leq \frac{1}{2}),$$

629 to obtain

$$\text{III} \leq \frac{8}{c(1 - e^{-(\theta - \frac{1}{2})c\eta_{\min,T}})} q.$$

630 Putting the three pieces together and reindexing the middle sum gives the first displayed bound
631 in the lemma. Putting the three pieces together and reindexing the middle sum gives the first
632 displayed bound in the lemma. The ‘‘in particular’’ version follows by $1 - e^{-x} \geq x/2$ with
633 $x = (\theta - \frac{1}{2})c\eta_{\min,T} \leq \frac{1}{2}$. \square

634 B.2 Proof of Main Result

635 *Proof of Theorem 3.*

636 From Lemma 17, for $c := \min\{\mu_{\text{PL}}/2, \rho/2\}$ and every t ,

$$\bar{\Xi}_{t+1} \leq \left(1 - \frac{c}{2}\eta_t\right)\bar{\Xi}_t + A\eta_t^2 + \frac{b^2}{2c}\eta_t + \frac{\tilde{H}}{n}\eta_t, \quad (20)$$

637 with the constants A, b, \tilde{H} defined in Lemma 17. Unrolling (20) and using $1 - \frac{c}{2}\eta_k \leq e^{-\frac{c}{2}\eta_k}$ yields

$$\bar{\Xi}_T \leq e^{-\frac{c}{2} \sum_{s=0}^{T-1} \eta_s} \bar{\Xi}_0 + A \sum_{t=0}^{T-1} \eta_t^2 e^{-\frac{c}{2} \sum_{s=t+1}^{T-1} \eta_s} + e^{\frac{c}{2} \eta_{\max, T}} \left(\frac{b^2}{c^2} + \frac{2\tilde{H}}{cn} \right), \quad (\star)$$

638 where $\eta_{\max, T} := \max_{0 \leq t < T} \eta_t$.

639 Lemma 15 (obtained from Lemma 13 (D_{weak})) gives

$$\bar{D}_T - \frac{C_{\text{hit}}}{n} \leq a \cdot \underbrace{\left(\sum_{t=0}^{T-1} \eta_t^2 e^{-2(2L - \frac{c}{2}) \sum_{s=t+1}^{T-1} \eta_s} \right)^{1/2}}_{\sqrt{S_u}} \cdot \underbrace{\left(\sum_{t=0}^{T-1} e^{-c \sum_{s=t+1}^{T-1} \eta_s} \mathbb{E}[\bar{\Xi}_t] \right)^{1/2}}_{\sqrt{S_v}} \quad (21)$$

640 where $C_{\text{hit}} := \frac{1}{\lambda} \left(2\tilde{B}_1 G + 2G \left(1 + \frac{L}{\rho} \right) \right)$. The potential sum, S_v , is bounded by Lemma 18. To
 641 combine the terms into the final compact form, we bound the products that arise after substitution.
 642 Since S_u is a convergent sum, it is bounded by a constant,

$$C_S := \sum_{t=0}^{\infty} \eta_t^2,$$

643 which depends only on the stepsize schedule $\{\eta_t\}$. Under the harmonic rule $\eta_t = \frac{c_1}{c_2 + t}$, $C_S =$
 644 $c_1^2 \sum_{t=0}^{\infty} (c_2 + t)^{-2} \leq c_1^2 \frac{\pi^2}{6}$, and for $c_2 > 1$ one may also use $C_S \leq \frac{c_1^2}{c_2 - 1}$. Moreover $S_u \rightarrow 0$ as
 645 $T \rightarrow \infty$, while S_v need not vanish; in the product we will use only the decaying initialization term
 646 and the variance component

$$S_v^{\text{var}}(T) := \sum_{t=0}^{T-1} \eta_t^2 e^{-\frac{c}{2} \sum_{s=t+1}^{T-1} \eta_s}.$$

647 This allows for two key bounds:

- 648 (i) The product involving the term $e^{-\frac{c}{2} \sum_{s=0}^{T-1} \eta_s} \bar{\Xi}_0$ is bounded by absorbing S_u into the constant:
 649 $S_u \cdot e^{-\frac{c}{2} \sum \eta_s} \Psi_{\alpha, 0}^{\max} \leq C_S \cdot e^{-\frac{c}{2} \sum \eta_s} \Psi_{\alpha, 0}^{\max}$.
- 650 (ii) The product of variance sums, $S_u \cdot S_v^{\text{var}}(T)$, is bounded by a multiple of their sum: $S_u \cdot$
 651 $S_v^{\text{var}}(T) \leq C_S \cdot S_v^{\text{var}}(T) \leq C_S \cdot (S_u + S_v^{\text{var}}(T))$.

652 Substituting these bounds, grouping all constants (a^2, C_S , etc.) into C_{var} , and relaxing the exponent⁸
 653 in S_u by defining $\kappa = \min\{\frac{3c}{4}, 2L - \frac{c}{2}\}$ yields the compact bound:

$$\begin{aligned} \bar{D}_T \leq & 2 C_{\text{dist}} \sqrt{e^{-\frac{c}{2} \sum_{s=0}^{T-1} \eta_s} \Psi_{\alpha, 0}^{\max} + C_{\text{var}} \left(L(1 + L/\rho) + \alpha \frac{L^2}{\rho} \right) G^2 \left[\sum_{t=0}^{T-1} \eta_t^2 e^{-2\kappa \sum_{s=t+1}^{T-1} \eta_s} + \sum_{t=0}^{T-1} \eta_t^2 e^{-\frac{c}{2} \sum_{s=t+1}^{T-1} \eta_s} \right]} \\ & + \frac{2G}{n} \left(\frac{(1 + L/\rho)^2}{\sqrt{\mu_{\text{PL}} \mu_{\text{QG}}}} + \frac{1}{\rho} \right) + \frac{C_{\text{hit}}}{n}, \end{aligned} \quad (22)$$

654 with $\kappa = \min\{\frac{3c}{4}, 2L - \frac{c}{2}\}$. □

⁸Since $x \mapsto e^{-\gamma x}$ is decreasing in $\gamma > 0$, replacing the larger decay parameter by the smaller $\kappa := \min\{2L - \theta, \frac{3}{4}c\}$ only increases the weighted sums, hence preserves a valid upper bound; we refer to this monotone weakening as “relaxing the exponent.”

655 *Proof of Corollary 4.* Let $\eta_t = \frac{c_1}{c_2+t}$ with $c_1 > 0$, $c_2 \geq 1$, and $c_1 < \min\{\frac{1}{2\kappa}, \frac{2}{c}\}$ so that the
 656 geometric constants below are finite. Then:

$$\sum_{s=0}^{T-1} \eta_s = c_1 \sum_{s=0}^{T-1} \frac{1}{c_2+s} \geq c_1 \log\left(\frac{c_2+T}{c_2}\right),$$

$$e^{-\frac{c}{2} \sum_{s=0}^{T-1} \eta_s} \leq \left(\frac{c_2}{c_2+T}\right)^{\frac{c c_1}{2}}, \quad e^{-\kappa \sum_{s=0}^{T-1} \eta_s} \leq \left(\frac{c_2}{c_2+T}\right)^{\kappa c_1},$$

$$\sum_{t=0}^{T-1} \eta_t^2 e^{-2\kappa \sum_{s=t+1}^{T-1} \eta_s} \leq \frac{c_1^2}{(1-2\kappa c_1)(c_2+T)}, \quad \sum_{t=0}^{T-1} \eta_t^2 e^{-\frac{c}{2} \sum_{s=t+1}^{T-1} \eta_s} \leq \frac{c_1^2}{(1-\frac{c}{2}c_1)(c_2+T)}.$$

657 Substituting these into the general bound above yields the explicit finite- T rate:

$$\begin{aligned} \varepsilon = \bar{D}_T \leq 2 C_{\text{dist}} & \left[\left(\frac{c_2}{c_2+T}\right)^{\frac{c_1 c}{2}} \Psi_{\alpha,0}^{\max} + \frac{C_{\text{var}} c_1^2 (L(1+L/\rho) + \alpha L^2/\rho) G^2}{(c_2+T)} \left(\frac{1}{1-2\kappa c_1} + \frac{1}{1-\frac{c}{2}c_1}\right) \right]^{1/2} \\ & + \frac{2G}{n} \left(\frac{(1+L/\rho)^2}{\sqrt{\mu_{\text{PL}} \mu_{\text{QG}}}} + \frac{1}{\rho}\right) + \frac{C_{\text{hit}}}{n}. \end{aligned} \quad (23)$$

658 Thus, for fixed n and feasible (c_1, c_2) , the optimization/stochastic term in (23) decays at rate $O((c_2 +$
 659 $T)^{-1/2})$, while the initialization bias decays polynomially as $(\frac{c_2}{c_2+T})^{\frac{c_1 c}{4}}$ inside the square root. \square

660 *Proof of Theorem 7.* Write $\hat{w} := \hat{w}^{(T)}$, $w_{\text{ERM}} := w^{\text{ERM}}$ and w^* for the population minimizer. Add
 661 and subtract the empirical risks of these three parameters:

$$\begin{aligned} \mathcal{R}(\hat{w}) - \mathcal{R}(w^*) &= \underbrace{[\mathcal{R}(\hat{w}) - \hat{\mathcal{R}}_n(\hat{w})]}_{\text{(A)}} + \underbrace{[\hat{\mathcal{R}}_n(\hat{w}) - \hat{\mathcal{R}}_n(w_{\text{ERM}})]}_{\text{(C)}} \\ &+ \underbrace{[\hat{\mathcal{R}}_n(w_{\text{ERM}}) - \hat{\mathcal{R}}_n(w^*)]}_{\text{(D)} \leq 0} + \underbrace{[\hat{\mathcal{R}}_n(w^*) - \mathcal{R}(w^*)]}_{\text{(B)}}. \end{aligned}$$

662 Theorem 6 gives

$$\mathbb{E}[(\text{A})] \leq (1+L/\rho)G\varepsilon_T.$$

663 Because w^* is deterministic, $\mathbb{E}[(\text{B})] = 0$. From Lemma 2, we have

$$(\text{C}) = \hat{\mathcal{R}}_n(\hat{w}) - \hat{\mathcal{R}}_n(w_{\text{ERM}}) \leq \frac{d_1}{d_2+T} \quad (\dagger)$$

664 (D) ≤ 0 deterministically (by definition of the ERM) and thus can be discarded when taking an upper
 665 bound.

666 Taking expectations and using $\mathbb{E}[(\text{B})] = 0$ and (\dagger) :

$$\mathbb{E}[\mathcal{R}(\hat{w}) - \mathcal{R}(w^*)] \leq (1+L/\rho)G\varepsilon_T + \frac{d_1}{d_2+T}.$$

667 \square

668 C Further discussions

669 C.1 Justification of Condition (C5)

670 **Lemma 19** (Bounded effective domain and bounded per-sample gradients for AGDA). *Fix a dataset*
 671 $\mathcal{D} = \{z_i\}_{i=1}^n$ *and a per-sample objective* $f(\cdot, \cdot; z)$. *Let the empirical saddle objective be*

$$F_{\mathcal{D}}(w, v) := \frac{1}{n} \sum_{i=1}^n f(w, v; z_i), \quad g_{\mathcal{D}}(w) := \max_v F_{\mathcal{D}}(w, v), \quad g_{\mathcal{D}}^* := \min_w g_{\mathcal{D}}(w).$$

672 Assume that $F_{\mathcal{D}}$ satisfies the following conditions (which are all shown to be satisfied in Kang et al.
673 [2025]), stated using the notation and constants of Yang et al. [2020]: there exist constants $l_Y > 0$
674 and $\mu_{1,Y}, \mu_{2,Y} > 0$ such that

675 (Y1) (Lipschitz gradient) $F_{\mathcal{D}}$ has l_Y -Lipschitz gradients in the joint variable (w, v) in the sense of
676 Yang et al. [2020, Assumption 1].

677 (Y2) (Existence of saddle point) $F_{\mathcal{D}}$ has at least one saddle point as in Yang et al. [2020, Assump-
678 tion 2].

679 (Y3) (Two-sided PL) $F_{\mathcal{D}}$ satisfies the two-sided PL condition with constants $\mu_{1,Y}, \mu_{2,Y}$ as in Yang
680 et al. [2020, Definition 2].

681 Run deterministic AGDA (i.e., Algorithm 1 of Yang et al. [2020] with $\sigma^2 = 0$) with constant step sizes

$$\tau_1 = \frac{\mu_{2,Y}^2}{18 l_Y^3}, \quad \tau_2 = \frac{1}{l_Y},$$

682 initialized at (w_0, v_0) :

$$w_{t+1} = w_t - \tau_1 \nabla_w F_{\mathcal{D}}(w_t, v_t), \quad v_{t+1} = v_t + \tau_2 \nabla_v F_{\mathcal{D}}(w_{t+1}, v_t).$$

683 Define the potential (same form as in Yang et al. [2020, Eq. (8)])

$$P_t := (g_{\mathcal{D}}(w_t) - g_{\mathcal{D}}^*) + \frac{1}{10} (g_{\mathcal{D}}(w_t) - F_{\mathcal{D}}(w_t, v_t)).$$

684 Then there exist a saddle point $(w_{\mathcal{D}}^*, v_{\mathcal{D}}^*)$ of $F_{\mathcal{D}}$ and a constant $\alpha_Y > 0$ (depending only on
685 $l_Y, \mu_{1,Y}, \mu_{2,Y}$) such that for all $t \geq 0$,

$$\|w_t - w_{\mathcal{D}}^*\|_2^2 + \|v_t - v_{\mathcal{D}}^*\|_2^2 \leq \alpha_Y \left(1 - \frac{\mu_{1,Y} \mu_{2,Y}^2}{36 l_Y^3}\right)^t P_0.$$

686 In particular, all iterates remain in the closed Euclidean ball

$$\Omega := \left\{ (w, v) : \|w - w_{\mathcal{D}}^*\|_2^2 + \|v - v_{\mathcal{D}}^*\|_2^2 \leq \alpha_Y P_0 \right\},$$

687 which is compact and convex.

688 Moreover, assume additionally that the sample space \mathcal{Z} is compact (e.g., finite) and that $(w, v, z) \mapsto$
689 $\nabla_w f(w, v; z)$ and $(w, v, z) \mapsto \nabla_v f(w, v; z)$ are continuous. Then the constant

$$G := \sup_{(w,v) \in \Omega, z \in \mathcal{Z}} \max \left\{ \|\nabla_w f(w, v; z)\|_2, \|\nabla_v f(w, v; z)\|_2 \right\}$$

690 is finite, and hence Condition (C5) in our paper holds for the AGDA iterates.

691 *Proof.* Under (Y1)–(Y3) and the specific choice of step sizes $\tau_1 = \mu_{2,Y}^2 / (18 l_Y^3)$ and $\tau_2 = 1 / l_Y$,
692 Theorem 3.2 of Yang et al. [2020] applies to the objective $F_{\mathcal{D}}$ and to the deterministic AGDA
693 recursion stated above. Therefore, there exist a saddle point $(w_{\mathcal{D}}^*, v_{\mathcal{D}}^*)$ of $F_{\mathcal{D}}$ and a constant $\alpha_Y > 0$
694 (depending only on $l_Y, \mu_{1,Y}, \mu_{2,Y}$) such that for all $t \geq 0$,

$$\|w_t - w_{\mathcal{D}}^*\|_2^2 + \|v_t - v_{\mathcal{D}}^*\|_2^2 \leq \alpha_Y \left(1 - \frac{\mu_{1,Y} \mu_{2,Y}^2}{36 l_Y^3}\right)^t P_0.$$

695 Since the contraction factor satisfies $0 < 1 - \frac{\mu_{1,Y} \mu_{2,Y}^2}{36 l_Y^3} < 1$, we have $\left(1 - \frac{\mu_{1,Y} \mu_{2,Y}^2}{36 l_Y^3}\right)^t \leq 1$ for all
696 $t \geq 0$. Hence

$$\|w_t - w_{\mathcal{D}}^*\|_2^2 + \|v_t - v_{\mathcal{D}}^*\|_2^2 \leq \alpha_Y P_0 \quad \forall t \geq 0.$$

697 This shows $(w_t, v_t) \in \Omega$ for all t . The set Ω is a closed Euclidean ball in a finite-dimensional space,
698 hence compact, and it is convex by construction.

699 By assumption, the maps $(w, v, z) \mapsto \nabla_w f(w, v; z)$ and $(w, v, z) \mapsto \nabla_v f(w, v; z)$ are continuous,
700 and $\Omega \times \mathcal{Z}$ is compact. Therefore, by the extreme value theorem, the function

$$(w, v, z) \mapsto \max \left\{ \|\nabla_w f(w, v; z)\|_2, \|\nabla_v f(w, v; z)\|_2 \right\}$$

701 attains its maximum on $\Omega \times \mathcal{Z}$, and the supremum

$$G = \sup_{(w,v) \in \Omega, z \in \mathcal{Z}} \max \left\{ \|\nabla_w f(w, v; z)\|_2, \|\nabla_v f(w, v; z)\|_2 \right\}$$

702 is finite. This is exactly Condition (C5) in our paper. \square

703 **C.2 Neural network parametrization and PL/QG condition.**

704 **Width/radius and high-probability regimes for uniform constants.** This subsection records the
 705 neural-network instantiation of Conditions (C1)–(C6), mainly (C2)–(C3). Fix a dataset $D = \{z_i\}_{i=1}^n$
 706 and write, as in Section 2.2,

$$F_D(w, v) := \frac{1}{n} \sum_{i=1}^n f(w, v; z_i), \quad \Phi_D(w) := \max_v F_D(w, v), \quad \Phi_D^* := \min_w \Phi_D(w).$$

707 Here w parametrizes the action-value network Q_w and v parametrizes the auxiliary function ζ_v in
 708 the bi-conjugate BRM objective (3). The PL/QG statements below are local: fix a failure level
 709 $\delta \in (0, 1)$ and a radius $R > 0$, and take the effective domain in Condition (C5) to be the initialization
 710 neighborhood

$$\Omega_R := \{(w, v) : \|(w, v) - (w_0, v_0)\| \leq R\} \cap (\mathcal{W} \times \mathcal{V}).$$

711 All constants used in the stability theorem are then understood on this same set: L is the joint
 712 smoothness constant of F_D , ρ is the strong concavity constant in v , G is the per-sample gradient
 713 bound, and $\mu_{\text{PL}}, \mu_{\text{QG}}$ are the PL and QG constants for Φ_D .

714 Consider a depth $(H + 1)$ fully connected network of width m , initialized with i.i.d. Gaussian weights.
 715 Let $\lambda_0 > 0$ denote the minimum eigenvalue of the empirical tangent kernel at initialization on the n
 716 sample inputs, and let $\rho_\sigma > 0$ be the lower bound on the derivative of the last-layer activation used in
 717 Liu et al. [2022]. Choose the target PL constant μ_{PL} so that

$$0 < \mu_{\text{PL}} < \rho_\sigma^2 \lambda_0.$$

718 Combining the over-parameterized neural-network PL result of Liu et al. [2022] with the BRM
 719 reduction of Kang et al. [2025], there are problem-dependent constants, with logarithmic factors
 720 hidden in $\tilde{\Omega}(\cdot)$, such that, if

$$m = \tilde{\Omega}\left(\frac{n R^{6H+2}}{(\lambda_0 - \mu_{\text{PL}}/\rho_\sigma^2)^2}\right),$$

721 then with probability at least $1 - \delta$ the value function Φ_D satisfies the PL inequality throughout the
 722 primal projection of Ω_R :

$$\frac{1}{2} \|\nabla \Phi_D(w)\|^2 \geq \mu_{\text{PL}} (\Phi_D(w) - \Phi_D^*).$$

723 Thus the constant appearing in Condition (C2) and in the contraction factor $c = \min\{\mu_{\text{PL}}/2, \rho/2\}$
 724 in Theorem 3 can be taken to be any admissible $\mu_{\text{PL}} \in (0, \rho_\sigma^2 \lambda_0)$ for which the displayed width
 725 condition is imposed.

726 We work on the intersection of this PL event with the event that the SGDA iterates remain inside Ω_R ;
 727 see Theorem 7 of Liu et al. [2022] for the corresponding radius/width regime for stochastic gradient
 728 methods. On this event, the softmax Bellman operator is smooth because $\eta > 0$, and the bi-conjugate
 729 BRM objective has a ρ -strongly concave inner maximization in the auxiliary block [Kang et al.,
 730 2025, Yang et al., 2020]. Since Φ_D is \mathcal{C}^2 on the local domain, the local PL–QG equivalence used in
 731 Section 2.2 yields a constant $\mu_{\text{QG}} > 0$ such that

$$\Phi_D(w) - \Phi_D^* \geq \frac{\mu_{\text{QG}}}{2} \|w - x_D^*\|^2$$

732 throughout the same domain, satisfying Condition (C3). For a neighboring dataset $D^{(i)}$, the same
 733 argument is applied to $F_{D^{(i)}}$ and $\Phi_{D^{(i)}}$; applying a union bound and then taking the larger upper-
 734 bound constants for L, G and the smaller lower-bound constants for $\rho, \mu_{\text{PL}}, \mu_{\text{QG}}$ gives the uniform
 735 constants required by Condition (S3).

736 For the rest of the conditions, under the same neural network Kang et al. [2025] proved that Conditions
 737 (C1) hold (Lemma 29), (C4) is relatively obvious since f in our BRM is simply a quadratic function
 738 of ζ with negative leading term, Yang et al. [2020] proved Condition (C6) holds for the problem of
 739 minimising Equation (5). Condition (C5) is in Appendix C.1.

740 **C.3 How index paring worked**

741 By synchronizing the minibatch selection, we neutralize it as a source of difference between the
742 two runs. Consequently, the parameter trajectories diverge only on the infrequent steps where the
743 replaced index i is sampled (a “hit”). On all other steps, the updates are identical, and the optimization
744 dynamics tend to pull the trajectories closer. The stability analysis thus bounds the cumulative effect
745 of these rare “hits” by balancing their small, infrequent perturbations against the constant, contractive
746 force of the optimization dynamics. This argument hinges entirely on the randomness of the sampling
747 process, which makes the “hits” probabilistic, and not on the statistical independence of the data
748 points, whose potential correlations are rendered irrelevant by the coupling.

749 **C.4 Uniformity setting**

750 The setting (S3) is standard in the stability literature: for example, Hardt et al. [2016] assumed each
751 per-example loss $f(\cdot; z)$ is L -Lipschitz and β -smooth uniformly in z , and Wang et al. [2022] assumed
752 the gradients and smoothness of $f(w, v; z)$ are bounded by global constants G and L for all z . These
753 conditions immediately imply that the corresponding constants are identical for any dataset.

754 **C.5 Discussions on contributions**

755 **What Kang et al. [2025] give us, and what they do not:** Kang et al. [2025] establish that, after a
756 bi-conjugate transformation, both the population and empirical BRM objectives enjoy PL-strongly-
757 concave structure in the parameterization of Q and ζ . They then prove global convergence of SGDA
758 to the empirical minimizer, and use PL to translate optimization error into parameter error. However,
759 they:

- 760 • do not analyze algorithmic stability,
- 761 • do not analyze generalization (population vs empirical BRM). In fact, in the latest version
762 of Kang et al., Lemma 28 references and uses our present result for the BRM sample
763 complexity. Thus, on their side, we only inherit the PL result and strong-concavity constants
764 for BRM; on our side, they use our sample-complexity result.

765 **What Wang et al. [2022] give us, and what they do not** Wang et al. [2022] develops on-average
766 argument stability for Markov chain SGD/SGDA and shows that for convex-concave or strongly-
767 convex-strongly-concave objectives, the excess population risk scales as $O(1/\sqrt{n})$, with explicit
768 dependence on the mixing parameter of the Markov chain over indices. We summarize the difference
769 between Wang et al. [2022] and our paper in Table 2. We reuse from Wang et al. [2022]:

- 770 • The concept of on-average argument stability and its connection to generalization (Lemma
771 5 in our notation).
- 772 • Some proof templates: two-run coupling on neighboring datasets and the idea of counting
773 "hits" of the replaced data point via an indicator $\mathbf{1}\{i_t = i\}$.

774 However, there are two crucial structural differences:

- 775 • Geometry: Wang et al. assume convex-concave or ρ -strongly-convex-strongly-concave
776 objectives. Our BRM objective is nonconvex in the primal parameters; only the value
777 function $\Phi_D(w) := \max_v F_D(w, v)$ is PL (and QG), and only in w . The algorithm,
778 however, is run on the bilevel saddle objective $F_D(w, v)$, not on Φ_D . This destroys the
779 standard convexity-based distance contraction used in Wang et al.
- Object of interest: Wang et al. work directly with the original minimax F and directly
analyze its risk/generalization. We instead are interested in the value function Φ_D (Bellman
residual) and the primal-dual gap built on Φ_D . The gradients used by SGDA at time t are
biased w.r.t. $\nabla\Phi_D(w_t)$ because the dual variable is not at its maximizer:

$$\Delta_t := \nabla_w F_D(w_t, v_t) - \nabla\Phi_D(w_t) \neq 0$$

780 This mismatch term Δ_t is absent in Wang et al. and is precisely what drives the need for our
781 Lyapunov potential.

Table 2: Comparison with prior work.

Work	Geometry	Setting	Stat. Rate	Opt. Rate
This paper	PL-strongly concave	Offline BRM	$\mathcal{O}(1/n)$	$\mathcal{O}((c_2 + T)^{-\min\{1/2, 3cc_1/8\}})$
Wang et al. (2022)	Convex-concave	minimax	$\mathcal{O}(1/\sqrt{n})$	$\mathcal{O}(1/\sqrt{T})$

782 C.6 Discussion on potential extensions

783 **Including experiments.** It is not a usual practice for theoretical papers discussing stability bounds
784 and corresponding generalization guarantees to include experiments, and therefore, not including
785 experiments is not a weakness of this paper. Please refer to Wang et al. [2022], Bousquet and Elisseeff
786 [2002], Charles and Papailiopoulos [2018], Feldman and Vondrak [2018].

787 **On non-smooth BRM.** Our analysis, which is based on *smooth* Bellman–residual objectives, does
788 not lose any generality for the entropy-regularized formulation adopted in Section 2. In particular,
789 when $\eta > 0$, the optimality equations involve the entropic log-sum-exp (softmax) operator

$$V_Q(s) = \eta \log \sum_{a \in \mathcal{A}} \exp(Q(s, a)/\eta),$$

790 so the Bellman operator and the resulting residual objective are differentiable and (under standard
791 boundedness conditions on the function class) admit Lipschitz gradients. By contrast, genuinely
792 non-smooth BRM variants (e.g., the hard max Bellman operator obtained in the zero-temperature
793 limit $\eta \rightarrow 0$, or objectives with non-differentiable losses/parameterizations) require different tools,
794 such as proximal or sub-gradient dynamics, smoothing/Moreau-envelope arguments, and stability
795 analyses tailored to non-smooth saddle problems, and are therefore orthogonal to the paper’s goal.

796 **Extending Bellman residual minimization to value function estimation.** Our main theorem
797 is stated for the Bellman Residual Minimization (BRM) objective. In offline RL, such a result is
798 often converted to a control-performance guarantee that is standard once one imposes a coverage
799 (a.k.a. concentrability) assumption that relates the state–action occupancy of the target/learned
800 policy to the data distribution. Under such coverage, a bound on the population Bellman residual
801 under the dataset distribution immediately yields a bound on the Bellman error along the policy’s
802 visitation distribution, which then translates to policy suboptimality via standard approximate dynamic
803 programming arguments. In fact, Kang et al. [2025] utilizes our result to instantiate this pipeline for
804 offline inverse reinforcement learning and related entropy-regularized control problems.

805 D Related Works

806 **Bellman Residual Minimization (BRM).** Bellman Residual Minimization (BRM) can be viewed
807 as a regression-style approach on offline reinforcement learning, in which the squared Bellman
808 residual plays the role of a regression loss. It has been identified as a particularly direct way to enforce
809 Bellman consistency under flexible function approximation, including nonparametric or neural-
810 network parameterizations [Jiang and Xie, 2024]. A central obstacle, however, is the double-sampling
811 problem [Antos et al., 2008]: the squared Bellman residual involves the square of a conditional
812 expectation, which cannot be unbiasedly estimated from a single next-state sample. An influential
813 remedy is to introduce a dual, or debiasing, correction through a bi-conjugate reformulation [Antos
814 et al., 2008, Dai et al., 2018, Patterson et al., 2022]. While this removes the statistical bias caused by
815 double sampling, it transforms BRM into a minimax optimization problem, whose primal component
816 is generally nonconvex under nonlinear function approximation. Consequently, BRM-style methods
817 have long been viewed as computationally delicate in offline RL, despite their conceptual appeal.

818 Recently, Kang et al. [2025] showed that, for entropy-regularized BRM and under neural-network ap-
819 proximation satisfying suitable over-parameterization conditions, the bi-conjugate objective admits a
820 PL–strongly-concave geometry and SGDA converges globally to an empirical saddle point. However,
821 their result is primarily an optimization guarantee: it does not directly establish algorithmic stability
822 or a statistical generalization bound comparing the empirical BRM objective with its population
823 counterpart. This statistical question is the main focus of the present work.

824 **Algorithmic stability and generalization.** Algorithmic stability is a classical approach to general-
825 ization analysis that controls how sensitively the output of a learning algorithm responds to small
826 perturbations of the training dataset [Bousquet and Elisseeff, 2002]. This framework has been devel-
827 oped through several notions of stability, including uniform stability, hypothesis stability, argument
828 stability, and on-average stability [Bousquet and Elisseeff, 2002, Hardt et al., 2016, Kuzborskij and
829 Lampert, 2018, Lei and Ying, 2020]. For gradient-based learning, Hardt et al. [2016] established
830 influential stability bounds for stochastic gradient descent, and subsequent works refined these bounds
831 using data-dependent or on-average notions of stability [Kuzborskij and Lampert, 2018, Lei and Ying,
832 2020, Lei et al., 2021]. Stability has also been used to study algorithms that converge to global optima
833 [Charles and Papailiopoulos, 2018] and to obtain refined high-probability generalization guarantees
834 for uniformly stable algorithms [Feldman and Vondrak, 2018, Bousquet et al., 2020].

835 More recently, stability-based generalization has been extended from minimization to minimax
836 optimization. This direction is particularly relevant for BRM because the bi-conjugate correction turns
837 the squared Bellman-residual objective into a saddle-point problem. Prior works studied the stability
838 and generalization of SGDA under various geometries, with recent extensions to decentralized,
839 distributed, and differentially private SGDA variants [Farnia and Ozdaglar, 2021, Lei et al., 2021,
840 Zhu et al., 2023, Zhang et al., 2024]. In parallel, Wang et al. [2022] developed an on-average
841 argument-stability framework for stochastic gradient methods with Markov-chain sampling and
842 established stability-to-generalization guarantees for both minimization and minimax problems. This
843 is especially relevant to offline RL, where data are often generated by trajectories and therefore need
844 not be i.i.d. However, these existing minimax stability results do not directly apply to the BRM
845 problem studied here. The key difficulty is geometric: the BRM saddle objective is not convex in
846 the primal parameter, and the algorithm is run on $F_D(w, v)$ rather than directly on the primal value
847 function $\Phi_D(w) := \max_v F_D(w, v)$. Thus the primal update contains a mismatch term

$$\nabla_w F_D(w_t, v_t) - \nabla \Phi_D(w_t),$$

848 unless the dual variable is already maximized. Our analysis addresses this mismatch through a
849 Lyapunov potential that couples the PL suboptimality of Φ_D with the dual gap, yielding an $O(1/n)$
850 on-average argument-stability bound and the corresponding BRM generalization guarantee.

851 **PL geometry and over-parameterized neural networks.** The Polyak–Łojasiewicz (PL) condition
852 is a central relaxation of strong convexity that still yields global convergence of gradient-based
853 methods despite nonconvexity [Polyak, 1963, Karimi et al., 2016]. It has become particularly
854 important in the analysis of over-parameterized models, where classical convexity-based arguments
855 are often inappropriate. In modern neural-network theory, over-parameterized loss landscapes have
856 been studied through several complementary perspectives, including the absence of strict spurious
857 local minima, connectivity of sublevel sets or global-minimum manifolds, neural tangent kernel
858 (NTK) dynamics, and local or global PL-type conditions [Du et al., 2019, Jacot et al., 2018, Lee et al.,
859 2019, Allen-Zhu et al., 2019, Arora et al., 2019, Soltanolkotabi et al., 2018, Liu et al., 2022, Chen
860 et al., 2023, Xu et al., 2025]. In particular, Liu et al. [2022] argue that sufficiently over-parameterized
861 nonlinear systems, including wide neural networks, satisfy a PL*-type condition in a neighborhood
862 of the initialization, and relate this property to tangent-kernel conditioning and Hessian norm control.
863 This line of work explains why gradient-based methods can converge globally in highly nonconvex
864 models without relying on convexity.

865 Our paper uses this perspective in the BRM setting: the relevant population and empirical soft
866 Bellman-residual objectives inherit a PL-strongly-concave structure under the neural-network
867 parametrization analyzed by Kang et al. [2025]. We then complement this optimization geome-
868 try with an algorithmic-stability analysis to obtain finite-sample statistical guarantees.

869 **Projected Mean-Squared Bellman Error (MSPBE)** Another approach is to minimize the pro-
870 jected mean-squared Bellman error (MSPBE), which first projects the Bellman residual onto the
871 approximation space and then measures its squared norm [Patterson et al., 2022]. This differs from
872 BRM, which directly minimizes the mean-squared Bellman error (MSBE) and therefore penalizes the
873 full Bellman inconsistency. MSPBE instead evaluates only the component of the Bellman residual
874 that lies within the representable function space. From an optimization perspective, this projection
875 can make the objective more tractable, since it only requires the residual to be small along directions
876 captured by the function class. However, this advantage comes at a cost: MSPBE ignores the compo-
877 nent of the Bellman residual orthogonal to the approximation space, and therefore does not control
878 the full Bellman inconsistency.

879 **Fitted Q-Iteration (FQI)** Fitted Q-Iteration (FQI), introduced by Ernst et al. [2005], is a canonical
880 value-based method for offline reinforcement learning. Although conceptually simple, iterative value-
881 based procedures are known to suffer from instability. In particular, approximate value iteration /
882 fitted value iteration can be unstable under off-policy sampling and function approximation, including
883 the case of infinite data and exact regression with linear function approximation. This instability
884 arises because FQI repeatedly fits value functions to targets that themselves depend on the current
885 estimate. Thus, the regression target evolves over iterations, creating an inherently unstable learning
886 process. More broadly, this phenomenon reflects the well-known deadly triad: the combination
887 of function approximation, bootstrapping, and off-policy learning [Tsitsiklis and Van Roy, 1996,
888 Van Hasselt et al., 2018, Sutton and Barto, 2018].

889 **Marginalized Importance Sampling (MIS)** Marginalized importance sampling (MIS) methods
890 cast offline reinforcement learning as a minimax problem by learning marginalized density ratios,
891 often interpreted as discriminators, to control worst-case reweighted Bellman errors [Nachum et al.,
892 2019, Zhang et al., 2020, Liu et al., 2018]. Unlike BRM, these methods optimize linear Bellman-error
893 objectives rather than squared residuals, and therefore avoid the double-sampling problem. However,
894 MIS objectives typically involve signed average Bellman errors. As a result, positive and negative
895 residuals can cancel each other out, potentially fail to penalize pointwise Bellman inconsistency
896 [Nachum et al., 2019, Zhang et al., 2020]. This distinction motivates our focus on BRM as a direct
897 residual-minimization approach, while addressing its statistical and optimization challenges through
898 the bi-conjugate formulation and stability analysis.