On the Convergence Rate of Density-Ratio Based Off-Policy Policy Gradient Methods

Jiawei Huang Department of Computer Science University of Illinois at Urbana-Champaign Urbana, IL 61801 jiaweih@illinois.edu Nan Jiang Department of Computer Science University of Illinois at Urbana-Champaign Urbana, IL 61801 nanjiang@illinois.edu

Abstract

We study the convergence properties of two optimization algorithms for off-policy policy gradient based on density-ratio learning. We establish general conditions that enable convergence and near-optimality guarantees, and show that these conditions can be satisfied in the linear case under standard assumptions. The keys to our analyses are the successful integration and application of stochastic first-order methods on solving saddle-point and non-convex optimization problems.

1 Introduction

Policy gradient (PG) is a very popular class of methods in empirical reinforcement-learning (RL) research, and has also attracted significant attention from the theoretical community recently [1]. Despite its appealing properties, classical PG typically requires on-policy roll-outs, making them not directly applicable to offline (or batch) RL. Recent development in marginalized importance sampling (MIS) methods [2, 3, 4, 5], however, has yielded promising off-policy policy-gradient estimators. For example, Nachum et al. [6] reformulated off-policy policy-optimization to a max-max-min problem, which faithfully optimizes the policy with sufficiently expressive function approximators [7]. A more general form of the problem considered by Yang et al. [5] is:

$$\max_{\pi \in \Pi} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}(\pi, w, Q) := \max_{\theta \in \Theta} \max_{\zeta \in Z} \min_{\xi \in \Xi} \mathcal{L}(\pi_{\theta}, w_{\zeta}, Q_{\xi})$$
$$:= (1 - \gamma) \mathbb{E}_{s_0 \sim \nu_0} [Q_{\xi}(s_0, \pi_{\theta})] + \mathbb{E}_{d^{\mu}} [w_{\zeta}(s, a) \Big(r + \gamma Q_{\xi}(s', \pi_{\theta}) - Q_{\xi}(s, a) \Big) + \lambda_O \mathbb{E}_{d^{\mu}} [f(Q_{\xi}(s, a))] - \lambda_w \mathbb{E}_{d^{\mu}} [g(w_{\zeta}(s, a))]$$
(1)

where π, w, Q are respectively parameterized by $(\theta, \zeta, \xi) \in \Theta \times Z \times \Xi$ (Θ, Z and Ξ are all convex sets), and we use Π, W, Q to denote their function classes; ν_0 is the initial state distribution, d^{μ} denotes the normalized discounted state-action occupancy induced by behavior policy μ (see Sec. 2.1

Despite the promising formulation, the problem takes a complex max-max-min form, which makes the optimization challenging. In this paper, we study the convergence guarantees of two natural optimization strategies for (the empirical version of) Eq.(2), and establish the conditions under which we can prove convergence rate and characterize the quality of the solutions. The actual objective, based on a sample D from d^{μ} , is

for a formal definition); $Q_{\xi}(s, \pi_{\theta})$ is short for $\mathbb{E}_{a \sim \pi_{\theta}(\cdot|s)}[Q_{\xi}(s, a)]$; f, g are regularizers.

$$\max_{\pi \in \Pi} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi, w, Q) := \max_{\theta \in \Theta} \max_{\zeta \in Z} \min_{\xi \in \Xi} \mathcal{L}^{D}(\pi_{\theta}, w_{\zeta}, Q_{\xi})$$
$$:= (1 - \gamma) \mathbb{E}_{s_{0} \sim \nu_{D}} [Q_{\xi}(s_{0}, \pi_{\theta})] + \mathbb{E}_{d^{D}} [w_{\zeta}(s, a) \Big(r + \gamma Q_{\xi}(s', \pi_{\theta}) - Q_{\xi}(s, a) \Big) \Big)$$
$$+ \frac{\lambda_{Q}}{2} \mathbb{E}_{d^{D}} [Q_{\xi}^{2}(s, a)] - \frac{\lambda_{w}}{2} \mathbb{E}_{d^{D}} [w_{\zeta}^{2}(s, a)].$$
(2)

Offline Reinforcement Learning Workshop at Neural Information Processing Systems, 2020.

Here we replace ν_0 with ν_D to denote the empirical initial distribution, and use d^D to denote the empirical state-action distribution in dataset. We also choose the regularizers to be quadratic functions.

In our analyses, we focus on the case when \mathcal{L}^D is strongly-concave w.r.t. ζ and strongly-convex w.r.t. ξ , but do not require the concavity related to θ . The strong concavity/convexity, among other assumptions we will introduce in Section 2.2, can be shown to be satisfied in the linear case under very standard assumptions (Appendix F).

Due to regularization, generalization error, and mis-specification error, there is inevitable bias between the stationary points of $\mathcal{L}^D(\pi_\theta, w_\zeta, Q_\xi)$ and $J(\pi_\theta)$, respectively, where $J(\pi_\theta)$ is the expected return of π_θ . Therefore, we focus on the convergence to the biased stationary point defined below.

Definition 1.1 (Biased stationary point).

$$\mathbb{E}[\|\nabla_{\theta} J(\pi_{\theta})\|] \le \varepsilon + \varepsilon_{data} + \varepsilon_{func} + \varepsilon_{reg} \tag{3}$$

where ε_{reg} , ε_{func} , ε_{data} are biases caused by regularization, mis-specified function class, and finitesample effects, respectively, as we will explain in Section 2. All norms in this paper is ℓ_2 norm unless specified otherwise. The expectation is over the randomness of the algorithm (e.g., the randomness in SGD) and not that of the data.

Paper Outline Our first algorithm, converts the original max-max-min problem to a max-min problem $\max_{(\theta,\zeta)\in\Theta\times Z} \min_{\xi\in\Xi} \mathcal{L}(\pi_{\theta}, w_{\zeta}, Q_{\xi})$, by simultaneously optimizing θ and ζ . Under the assumptions identified in Section 2.2, we prove that the stationary point returned by any stochastic optimization algorithm for non-convex-strongly-concave problems is also a biased stationary point in Definition 1.1. As a result, the $O(\varepsilon^{-3})$ convergence rate can be established based on a recent result on non-convex-strongly-concave optimization [8].

We then study another algorithm, where we iteratively solve the inner strongly-concave-stronglyconvex max-min problem $\max_{\zeta \in \mathbb{Z}} \min_{\xi \in \Xi} \mathcal{L}(\pi_{\theta}, w_{\zeta}, Q_{\xi})$ for fixed θ and the outer non-convex optimization problem $\max_{\theta \in \Theta} \mathcal{L}(\pi_{\theta}, w_{\zeta}, Q_{\xi})$ for fixed ζ and ξ . For the inner loop, we assume an oracle that solves the saddle-point problem, and provide a concrete example in Appendix E. For the outer loop, the main technique difficulty is that, the loss function $\mathcal{L}(\pi_{\theta}, w_{\zeta_t}, Q_{\xi_t})$ varies across iterations because we update ζ_t, ξ_t in the inner loop, which prevents us from adapting existing non-convex optimization algorithms directly. We resolve this difficulty by coordinating the inner and the outer loops so that we can relate the variation $\|\zeta_{t+1} - \zeta_t\|$ and $\|\xi_{t+1} - \xi_t\|$ with $\|\theta_{t+1} - \theta_t\|$. The convergence rate to a biased stationary point of our second strategy is $O(\varepsilon^{-4})$.

1.1 Related works

Recently, there has been a lot of interest in turning MIS methods for off-policy evaluation [3, 9, 2] into off-policy policy-optimization algorithms. Liu et al. [10] presented OPPOSD with convergence guarantees, but the convergence relies on accurately estimating the density ratio and the value function via MIS, which were treated as a black box without further analysis. [6, 7] discussed policy optimization given arbitrary off-policy dataset, but no convergence analysis was performed. Another style of off-policy policy-improvement algorithms is off-policy actor-critic [11, 12, 13]. Although [13] presented a provably convergent algorithm, where only asymptotic convergence was proved and no finite convergence rate was given.

Meanwhile, along with the progress of the variance reduction techniques for non-convex optimization, there are several emerging works analyzing convergence rates in RL settings [14, 15, 16, 17, 18]. However, all of them require on-policy interaction with the environment, whereas our focus is the off-policy setting.

2 Preliminary

2.1 Markov Decision Process

We consider an infinite-horizon discounted MDP $(S, A, R, P, \gamma, \nu_0)$, where S and A are the state and action spaces, respectively, which we assume to be finite but can be arbitrarily large. $R: S \times A \rightarrow \Delta([0, 1])$ is the reward function. $P: S \times A \rightarrow \Delta(S)$ is the transition function, γ is the discount factor and ν_0 denotes the initial state distribution.

For arbitrary policy π , we use $d^{\pi}(s,a) = (1-\gamma)\mathbb{E}_{\tau \sim \pi, s_0 \sim \nu_0}[\sum_{t=0}^{\infty} \gamma^t p(s_t = s, a_t = a)]$ to denote the normalized discounted state-action occupancy, where $\tau \sim \pi, s_0 \sim \nu_0$ means a trajectory $\tau = \{s_0, a_0, s_1, a_1, ...\}$ is sampled according to the rule that $s_0 \sim \nu_0, a_0 \sim \pi(\cdot|s_0), s_1 \sim P(\cdot|s_0, a_0), a_1 \sim \pi(\cdot|s_1), ...,$ and $p(s_t = s, a_t = a)$ denotes the probability that the *t*-th state-action pair are exactly (s, a). We also use $Q^{\pi}(s, a) = \mathbb{E}_{\tau \sim \pi, s_0 = s, a_0 = a}[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$ to denote the Q-function of π . It is well-known that Q^{π} satisfies the Bellman Equation:

$$Q^{\pi}(s,a) = \mathcal{T}^{\pi}Q^{\pi}(s,a) := \mathbb{E}_{r \sim R(s,a), s' \sim P(\cdot|s,a), a' \sim \pi(\cdot|s')} [r + \gamma Q^{\pi}(s',a')].$$

Define $J(\pi) = \mathbb{E}_{s \sim \nu_0, a \sim \pi(\cdot|s_0)}[Q^{\pi}(s, a)] = \frac{1}{1-\gamma}\mathbb{E}_{s, a \sim d^{\pi}}[r(s, a)]$ as the expected return of policy π . If π is parameterized by θ and differentiable, the policy-gradient theorem [19] states that

$$\nabla_{\theta} J(\pi_{\theta}) = \frac{1}{1 - \gamma} \mathbb{E}_{s, a \sim d^{\pi}} [Q^{\pi}(s, a) \nabla_{\theta} \log \pi(a|s)].$$

In the off-policy setting, we can only get access to d^{μ} , the discounted state-action occupancy w.r.t. another policy μ . Then we can rewrite $\nabla_{\theta} J(\pi)$ by introducing the importance ratio $w^{\pi}(s, a) := \frac{d^{\pi}(s, a)}{d^{\mu}(s, a)}$.

$$\nabla_{\theta} J(\pi_{\theta}) = \frac{1}{1-\gamma} \mathbb{E}_{s,a \sim d^{\mu}} [w^{\pi}(s,a) Q^{\pi}(s,a) \nabla_{\theta} \log \pi(a|s)].$$

In the rest of the paper, we will refer μ as the behavior policy, and refer π as the target policy whose performance we are interested in.

In practice, usually, we are only provided with an off-line dataset instead of the exact distribution d^{μ} , which we denote as $D = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^{|D|}$. Each tuple is sampled by $s_i, a_i \sim d^{\mu}, r_i \sim R(s_i, a_i), s'_i \sim P(\cdot|s_i, a_i)$, and we use d^D to denote the empirical state-action distribution.

2.2 Assumptions and Definitions

We now introduce the assumptions and definitions that will later enable us to establish the convergence guarantees and characterize the solution quality. We will also introduce some algorithm-specific assumptions later. While some of the assumptions (e.g., Assumption C) are quite strong, in Appendix F we show they are automatically satisfied in the linear setting under more standard assumptions.

Assumption A (Smoothness).

(a) For any $s, a \in S \times A$ and $\theta \in \Theta$, $\pi_{\theta}(s, a)$ is second-order differentiable w.r.t. θ , and there exist constants G and H, s.t.

$$\|\nabla_{\theta} \log \pi_{\theta}(a|s)\| \le G, \qquad \|\nabla_{\theta}^{2} \log \pi_{\theta}(a|s)\|_{op} \le H$$
(4)

where $\|\cdot\|_{op}$ is the matrix operator norm.

(b) For any $\xi, \xi_1, \xi_2 \in \Xi, \zeta, \zeta_1, \zeta_2 \in Z, (s, a) \in S \times A$, there are constants C_Q, C_W, L_Q, L_w , s.t. $|Q_{\xi}(s, a)| \leq C_Q; \quad |Q_{\xi_1}(s, a) - Q_{\xi_2}(s, a)| \leq L_Q ||\xi_1 - \xi_2||;$ $|w_{\zeta}(s, a)| \leq C_W; \quad |w_{\zeta_1}(s, a) - w_{\zeta_2}(s, a)| \leq L_w ||\zeta_1 - \zeta_2||;$

Usually, in practice, we normalize the expectation of w_{ζ} to 1, so $C_{W} > 1$ in general.

(c) Let $v \in V = \Theta \times Z \times \Xi$ denote a vector formed by concatenating θ, ζ, ξ . For any $v, v_1, v_2 \in V$, \mathcal{L}^D defined in Eq.(2) is differentiable w.r.t. v, and there exists constant L s.t.

$$\begin{aligned} \|\nabla_{v}\mathcal{L}^{D}(v_{1}) - \nabla_{v}\mathcal{L}^{D}(v_{2})\| : \\ = \|\nabla_{\theta}\mathcal{L}^{D}(v_{1}) - \nabla_{\theta}\mathcal{L}^{D}(v_{2})\| + \|\nabla_{\zeta}\mathcal{L}^{D}(v_{1}) - \nabla_{\zeta}\mathcal{L}^{D}(v_{2})\| + \|\nabla_{\xi}\mathcal{L}^{D}(v_{1}) - \nabla_{\xi}\mathcal{L}^{D}(v_{2})\| \\ \le L\|\theta_{1} - \theta_{2}\| + L\|\zeta_{1} - \zeta_{2}\| + L\|\xi_{1} - \xi_{2}\| \end{aligned}$$

Assumption B (Exploratory Data). Recall the behavior policy is denoted as μ . We assume there exists a constant C > 0, for arbitrary $\pi \in \Pi$ and any $(s, a) \in S \times A$, we have

$$w^{\pi}(s,a) := \frac{d^{\pi}(s,a)}{d^{\mu}(s,a)} \le C, \qquad w^{\pi}_{d^{\mu}}(s,a) := \frac{d^{\pi}_{d^{\mu}}(s,a)}{d^{\mu}(s,a)} \le C$$

where $d_{d\mu}^{\pi}(s,a) := (1-\gamma)\mathbb{E}_{\tau \sim \pi, s_0, a_0 \sim d^{\pi}(\cdot, \cdot)}[\sum_{t=0}^{\infty} \gamma^t p(s_t = s, a_t = a)]$ is the normalized discounted state-action occupancy by treating d^{μ} as initial distribution.

Assumption C (Strongly-Convex-Strongly-Concave). We use u_Z and u_{Ξ} to denote the dimension of vector parameters ζ and ξ . Given arbitrary $\theta \in \Theta, \zeta \in \mathbb{R}^{u_Z}, \mathcal{L}^D(\theta, \zeta, \cdot)$ is μ_{ξ} -strongly convex w.r.t. $\xi \in \Xi$. Given arbitrary $\theta \in \Theta, \xi \in \mathbb{R}^{u_{\Xi}}, \mathcal{L}^D(\theta, \cdot, \xi)$ is μ_{ζ} -strongly concave w.r.t. $\zeta \in Z$.

Remark 2.1. In fact, the regularization terms is necessary if we want Assumption C to hold when one of w^{π} and Q^{π} is realizable. We defer the discussion to Appendix B.

Assumption D. Denote $(\zeta_{\theta}^*, \xi_{\theta}^*)$ as the saddle point of $\mathcal{L}^D(\theta, \zeta, \xi)$ without constraint on ζ and ξ . For arbitrary π_{θ} parameterized by $\theta \in \Theta$, $(\zeta_{\theta}^*, \xi_{\theta}^*) \in Z \times \Xi$.

Remark 2.2. Based on Assumption A, C, since both Z and Ξ are convex sets, Assumption D implies that

$$\|\nabla_{\zeta} \mathcal{L}^D(\theta, \zeta_{\theta}^*, \xi_{\theta}^*)\| = \|\nabla_{\xi} \mathcal{L}^D(\theta, \zeta_{\theta}^*, \xi_{\theta}^*)\| = 0$$

Definition 2.3 (Generalization Error). Suppose there exists a constant ε'_{data} , for arbitrary $\pi_{\theta}, w_{\zeta}, Q_{\xi} \in \Pi \times \mathcal{W} \times \mathcal{Q}$, we have:

$$\begin{aligned} |\mathcal{L}(\pi_{\theta}, w_{\zeta}, Q_{\xi}) - \mathcal{L}^{D}(\pi_{\theta}, w_{\zeta}, Q_{\xi})| &\leq \varepsilon'_{data} \\ \|\nabla_{\theta} \mathcal{L}(\pi_{\theta}, w_{\mu}^{*}, Q_{\mu}^{*}) - \nabla_{\theta} \mathcal{L}^{D}(\pi_{\theta}, w_{\mu}^{*}, Q_{\mu}^{*})\|^{2} &\leq \varepsilon'_{data} \end{aligned}$$

where $(w_{\mu}^*, Q_{\mu}^*) := \arg \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}(\pi, w, Q).$

Proposition 2.4. Denote $\varepsilon_{data} := (2\kappa_{\zeta}\kappa_{\xi} + 2\kappa_{\zeta} + 2\kappa_{\zeta} + \sqrt{2}/2)\sqrt{2\varepsilon'_{data}}$, where ε'_{data} is defined in Definition 2.3, $\kappa_{\zeta} = L/\mu_{\zeta}$, $\kappa_{\xi} = L/\mu_{\xi}$. Under Assumption A and C, we have:

$$\left\|\nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}(\pi_{\theta}, w, Q) - \nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi_{\theta}, w, Q)\right\| \leq \varepsilon_{data}$$

We defer the proof to Appendix A.

Definition 2.5 (Mis-specification Error).

(1) For arbitrary $\pi \in \Pi$, denote $w_{\zeta^{\pi}} := \arg \min_{w \in \mathcal{W}} \|w - w_{\mathcal{L}}^{\pi}\|_{\Lambda}^2$ parameterized by $\zeta^{\pi} \in Z$, where $w_{\mathcal{L}}^{\pi} = \arg \max_{w \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}} \min_{Q \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}} \mathcal{L}(\pi, w, Q)$. We define

$$\varepsilon_1 := \max_{\pi \in \Pi} \| w_{\zeta^{\pi}} - w_{\mathcal{L}}^{\pi} \|_{\Lambda}^2$$

(2) For arbitrary policy $\pi \in \Pi$ and $w \in W$, denote $Q_{\xi_w^{\pi}} := \arg \min_{Q \in \mathcal{Q}} \mathcal{L}(\pi, w, Q)$ parameterized by $\xi_w^{\pi} \in \Xi$. We define

$$\varepsilon_2 := \max_{w \in \mathcal{W}, \pi \in \Pi} \|Q_{\xi_w^{\pi}} - \arg \min_{Q \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}} \mathcal{L}(\pi, w, Q)\|_{\Lambda}^2$$

A consequence of Assumptions A and C is Proposition 2.6, that we can use ε_1 and ε_2 defined in Definition 2.5 to bound the weighted difference between the saddle points of $\mathcal{L}^D(\pi, w, Q)$ with and without constraining w and Q on $\mathcal{W} \times Q$, respectively, which is crucial to analyzing the bias resulting from the mis-specified function classes. We defer its proof to Appendix A.

Proposition 2.6. Under Assumption A and C, for arbitrary $\pi \in \Pi$, we have:

$$\mathbb{E}_{d^{\mu}}[|w_{\mu}^{*}(s,a) - w_{\mathcal{L}}^{\pi}(s,a)|^{2}] \leq \varepsilon_{\mathcal{W}} := 4\frac{\lambda_{\max}^{2}}{\lambda_{Q}\lambda_{w}}\varepsilon_{1} + 2\frac{L_{w}^{2}\lambda_{\max}}{\mu_{\zeta}}\varepsilon_{2}$$
$$\mathbb{E}_{d^{\mu}}[|Q_{\mu}^{*}(s,a) - Q_{\mathcal{L}}^{\pi}(s,a)|^{2}] \leq \varepsilon_{\mathcal{Q}} := 8\frac{\lambda_{\max}^{3}}{\lambda_{Q}^{2}\lambda_{w}}\varepsilon_{1} + (2 + 4\frac{L_{w}^{2}\lambda_{\max}^{2}}{\lambda_{Q}\mu_{\zeta}})\varepsilon_{2}$$

where $(w_{\mu}^{*}, Q_{\mu}^{*})$ denotes the saddle point of $\mathcal{L}(\pi, w, Q)$ constrained by $w, Q \in \mathcal{W} \times \mathcal{Q}$, $(w_{\mathcal{L}}^{\pi}, Q_{\mathcal{L}}^{\pi})$ denotes the saddle point of $\mathcal{L}(\pi, w, Q)$ without any constraint on w and Q, $\lambda_{\max} = \max\{\lambda_Q, \lambda_w\}$, L_w is defined in Assumption A, μ_{ζ} is defined in Assumption C.

2.3 Main goal of the analyses

First, by applying the triangle inequality, we have:

$$\|\nabla_{\theta} J(\pi_{\theta})\| \leq \|\nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi_{\theta}, w, Q)\| + \|\nabla_{\theta} J(\pi_{\theta}) - \nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi_{\theta}, w, Q)\|$$

where w^*, Q^* denotes the saddle point of $\mathcal{L}^D(\pi_\theta, w, Q)$ constrained by $w, Q \in \mathcal{W} \times \mathcal{Q}$. Optimizing the loss function $\mathcal{L}^D(\pi, w, Q)$ may offer us a better θ to decrease the first term, while based on above Assumptions, we can bound the second term in the following Theorem.

Theorem 2.7. [Bias] Under Assumption A, B, C, given arbitrary $\theta \in \Theta$, we have

$$\|\nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi_{\theta}, w, Q) - \nabla_{\theta} J(\pi_{\theta})\| \leq \varepsilon_{reg} + \varepsilon_{func} + \varepsilon_{data}$$

where ε_{data} is defined in Proposition 2.4, and

$$\varepsilon_{func} = \frac{G}{1 - \gamma} \Big(\sqrt{C\varepsilon_{\mathcal{Q}}} + C_{\mathcal{W}} \sqrt{\frac{\gamma \varepsilon_{\mathcal{Q}} C}{1 - \gamma}} + \sqrt{\frac{\gamma \varepsilon_{\mathcal{Q}} \varepsilon_{\mathcal{W}} C}{1 - \gamma}} + \gamma C_{\mathcal{Q}} \sqrt{\varepsilon_{\mathcal{W}}} \Big)$$

$$(\varepsilon_{\mathcal{W}} \text{ and } \varepsilon_{\mathcal{Q}} \text{ defined in Prop. 2.6})$$

We defer its proof to Appendix B.

As we can see, $\|\nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi_{\theta}, w, Q) - \nabla_{\theta} J(\pi_{\theta})\|$ can be controlled by three terms. ε_{data} reflects the generalization error, and should be small if we have plenty of data. ε_{reg} depends on the magnitude of regularization, and will decrease as λ_w and λ_Q . As for ε_{func} , it depends on the approximation error $\varepsilon_{\mathcal{W}}$ and $\varepsilon_{\mathcal{Q}}$, which are proportional to ε_1 and ε_2 . Besides, because μ_{ζ} should be proportional to λ_w and L_w does not depend on regularization, the coefficients before ε_1 and ε_2 should not vary a lot as we change λ_w and λ_Q while keeping $\lambda_w \approx \lambda_Q$ (but ε_1 and ε_2 may change with λ_w and λ_Q). In general, a larger dataset, better function classes and smaller λ_w and λ_Q may result in smaller bias, while smaller regularization can lead to weaker strong-concavity or strong-convexity of the loss function and make the convergence slower.

Based on the discussion above, our goal is to find stochastic optimization algorithms, which can return us π_{θ} after consuming $Poly(\varepsilon^{-1})$ samples from dataset (we omit the dependence on others such as μ_{ζ}, μ_{ξ} and etc.), satisfying the following biased stationary condition in Definition 1.1:

$$\mathbb{E}[\|\nabla_{\theta} J(\pi_{\theta})\|] \le \varepsilon + \varepsilon_{data} + \varepsilon_{func} + \varepsilon_{reg}$$
(5)

where ε_{data} is defined in 2.3 and ε_{func} and ε_{reg} are defined in Theorem 2.7.

Since D can be extremely large, we consider stochastic optimization, and introduce another crucial assumption about the stochastic gradient:

Assumption E (Variance of Estimated Gradient). We use $\mathbb{E}_{s,a,r,s',a_0,a'}[\cdot]$ as a short note of

$$\mathbb{E}_{(s,a,r,s')\sim d^D,a_0\sim\pi(\cdot|s),a'\sim\pi(\cdot|s')}[\cdot]$$

and use $\mathcal{L}^{(s,a,r,s',a_0,a')}(\theta,\zeta,\xi)$ to denote the gradient estimation with only one sample defined by:

$$(1-\gamma)Q_{\xi}(s,a_{0})\pi_{\theta}(a_{0}|s)\mathbb{I}[s\in S_{0}] + w_{\zeta}(s,a)\Big(r+\gamma Q_{\xi}(s',a')\pi_{\theta}(a'|s') - Q_{\xi}(s,a)\Big) + \frac{\lambda_{Q}}{2}Q_{\xi}^{2}(s,a) - \frac{\lambda_{w}}{2}w_{\zeta}^{2}(s,a)$$

where $\mathbb{I}[s \in S_0]$ equals 1 only if s is generated at the first step in a trajectory and equals 0 otherwise (note that we allow the case when a state in the initial state sets can be visited at step $t \ge 1$.). We assume that, there exists a positive constant σ , for arbitrary $\theta, \zeta, \xi \in \Theta \times Z \times \Xi$, we have:

$$\begin{split} & \mathbb{E}_{s,a,r,s',a_0,a'}[\|\nabla_{\theta}\mathcal{L}^{(s,a,r,s',a_0,a')}(\theta,\zeta,\xi) - \nabla_{\theta}\mathcal{L}^D(\theta,\zeta,\xi)\|^2] \leq \sigma^2 \\ & \mathbb{E}_{s,a,r,s',a_0,a'}[\|\nabla_{\zeta}\mathcal{L}^{(s,a,r,s',a_0,a')}(\theta,\zeta,\xi) - \nabla_{\zeta}\mathcal{L}^D(\theta,\zeta,\xi)\|^2] \leq \sigma^2 \\ & \mathbb{E}_{s,a,r,s',a_0,a'}[\|\nabla_{\xi}\mathcal{L}^{(s,a,r,s',a_0,a')}(\theta,\zeta,\xi) - \nabla_{\xi}\mathcal{L}^D(\theta,\zeta,\xi)\|^2] \leq \sigma^2 \end{split}$$

Remark 2.8. The upper bound on the variance of the gradients w.r.t. θ , ζ and ξ are usually assumed to be different. Here we use σ to refer to the maximum of these upper bounds to simplify notations.

3 Strategy 1: Converting Max-Max-Min to Max-min problem

A heuristic optimization strategy for (2) is to rewrite the original max-max-min problem $\max_{\theta} \max_{\zeta} \min_{\xi} \mathcal{L}^{D}(\theta, \zeta, \xi)$ to a max-min problem $\max_{\theta, \zeta} \min_{\xi} \mathcal{L}^{D}(\theta, \zeta, \xi)$. Given Assumption A and C, we know $\max_{\theta, \zeta} \min_{\xi} \mathcal{L}^{D}(\theta, \zeta, \xi)$ is a standard non-concave-strongly-convex problem, which can be solved efficiently based on the recent progress on non-convex-strongly-concave optimization [20, 8].

In this section, we prove the equivalence between the stationary point of the non-convex-stronglyconcave saddle-point problem and the stationary point of our policy gradient objective:

Theorem 3.1. [Equivalence Between Stationary Points] Under Assumption A, C and D, suppose there exists a $\theta \in \Theta$ s.t. $\|\nabla_{\theta} \max_{\zeta \in Z} \min_{\xi \in \Xi} \mathcal{L}^{D}(\theta, \zeta, \xi)\| = 0$ and there is an Algorithm provides us one stationary point $(\theta_T, \zeta_T, \xi_T)$ of the non-concave-strongly-convex problem $\max_{\theta, \zeta} \min_{\xi} \mathcal{L}^{D}(\theta, \zeta, \xi)$ after running T iterations, which statisfying the following conditions in expectation over the randomness of algorithm.

$$\mathbb{E}[\|\nabla_{\theta,\zeta}\mathcal{L}^{D}(\theta_{T},\zeta_{T},\phi_{\theta_{T}}(\zeta_{T}))\|]$$

:= $\mathbb{E}[\|\nabla_{\theta}\mathcal{L}^{D}(\theta_{T},\zeta_{T},\phi_{\theta_{T}}(\zeta_{T}))\| + \|\nabla_{\zeta}\mathcal{L}^{D}(\theta_{T},\zeta_{T},\phi_{\theta_{T}}(\zeta_{T}))\|] \le \frac{\varepsilon}{(\kappa_{\xi}+1)(\kappa_{\zeta}+1)}$ (6)

where $\phi_{\theta}(\zeta) = \arg \min_{\xi \in \Xi} \mathcal{L}^{D}(\theta, \zeta, \xi)$. Then, we have

$$\mathbb{E}[\|\nabla_{\theta} J(\pi_{\theta_T})\|] \le \varepsilon + \varepsilon_{data} + \varepsilon_{func} + \varepsilon_{reg}$$

In Appendix C, we will give the detailed proof. Besides, we also list algorithm examples which can return us stationary points satisfying Eq.(6).

4 Strategy 2: Stochastic Recursive Momentum with Saddle-Point Oracle

In this section, we propose a new algorithm, based on stochastic recursive momentum and a saddlepoint oracle. We will provide a concrete example of the oracle algorithm in the Appendix E.

Definition 4.1 (Oracle Algorithm). Suppose we have an oracle algorithm *Oracle*. For arbitrary strongly-concave-strongly-convex problem $f(\zeta, \xi)$ with saddle point $(\zeta^*, \xi^*) \in Z \times \Xi$, and arbitrary $0 < \beta \leq 1$ and c > 0, starting from a random initializer $(\zeta_0, \xi_0) \in Z \times \Xi$ and executing finite steps, *Oracle* returns a solution (ζ_K, ξ_K) satisfying

$$\mathbb{E}[\|\zeta_K - \zeta^*\|^2 + \|\xi_K - \xi^*\|^2] \le \frac{\beta}{2} \mathbb{E}[\|\zeta_0 - \zeta^*\|^2 + \|\xi_0 - \xi^*\|^2] + c$$
(7)

Next, we present our oracle based stochastic recursive momentum algorithm (O-SRM), inspired by the on-policy SRM [17]. In our algorithm, we choose $\Theta = \mathbb{R}^{u_{\Theta}}$ where u_{Θ} is the dimension of Θ . As a result, we will not do projection after update θ and there must exist stationary points of $J(\pi_{\theta})$ and $\max_{\zeta \in \mathbb{Z}} \min_{\xi \in \Xi} \mathcal{L}^D(\theta, \zeta, \xi)$. We will use $\nabla_{\theta} \mathcal{L}^B(\theta, \zeta, \xi)$ as a short note of the empirical version of the gradient estimator, i.e.

$$\nabla_{\theta} \mathcal{L}^{B}(\theta,\zeta,\xi) = \frac{1}{|B|} \sum_{B} (1-\gamma)Q(s^{i},a_{0}^{i})\nabla_{\theta}\log\pi(a_{0}^{i}|s^{i})\mathbb{I}[s^{i}\in S_{0}] + \gamma w(s^{i},a^{i})Q(s^{\prime i},a^{\prime i})\nabla_{\theta}\log\pi(a^{\prime i}|s^{\prime i})$$

where (s^i, a^i, r^i, s'^i) for i = 1, 2, ..., |B| are elements in B sampled from d^D , and $a_0^i \sim \pi(\cdot |s^i), a'^i \sim \pi(\cdot |s'^i)$.

Algorithm 1: O-SRM

1 Input: Total number of iteration T; Learning rate $\eta_{\theta}, \eta_{\zeta}, \eta_{\xi}$; Dataset distribution d^{D} ; Oracle parameter β . 2 Initialize $\theta_{0}, \zeta_{-1}, \xi_{-1}$ 3 $\zeta_{0}, \xi_{0} \leftarrow \operatorname{Oracle}(T_{1}, \eta_{\zeta}, \eta_{\xi}, \theta_{0}, \zeta_{-1}, \xi_{-1}, d^{D})$ 4 Sample $B_{0} \sim d^{D}$ with batch size $|B_{0}|$ and estimate $g_{\theta}^{0} = \nabla_{\theta} \mathcal{L}^{B_{0}}(\theta_{0}, \zeta_{0}, \xi_{0})$ 5 for t = 0, 1, 2, ...T - 1 do 6 $| \theta_{t+1} \leftarrow \theta_{t} + \eta_{\theta} g_{\theta}^{t}$ 7 $| \zeta_{t+1}, \xi_{t+1} \leftarrow Oracle(\beta, \theta_{t+1}, \zeta_{t}, \xi_{t}, d^{D}, \beta)$ 8 Sample $B \sim d^{D}$; 9 $| g_{\theta}^{t+1} = (1 - \alpha) \left(g_{\theta}^{t} - \nabla_{\theta} \mathcal{L}^{B}(\theta_{t}, \zeta_{t}, \xi_{t}) \right) + \nabla_{\theta} \mathcal{L}^{B}(\theta_{t+1}, \zeta_{t+1}, \xi_{t+1})$ 10 end 11 Output: Sample $\theta_{out} \sim \operatorname{Unif}\{\theta_{0}, \theta_{1}, ..., \theta_{T}\}$ and output π_{θ} .

4.1 Additional Assumptions for Algorithm 1

Assumption F (Diameter). We use Z and Ξ to denote the sets of parameters ζ and ξ , respectively, we assume Z and Ξ are both convex and bounded set, and there exists a constant d, such that the diameters of Z and Ξ are bounded by d.

4.2 Algorithm Analysis

We first derive the smoothness of $J(\pi_{\theta})$:

Proposition 4.2. Under Assumption A, $J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}, s_0 \sim \nu_0} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$ is L_J smooth with

$$L_J := \frac{H}{(1-\gamma)^2} + \frac{(1+\gamma)G^2}{(1-\gamma)^3}$$

Theorem 4.3. Under Assumption A-F and H, given arbitrary ε , by choosing Algorithm 3 as the Oracle, Algorithm 1 will return us a policy $\pi_{\theta_{out}}$, satisfying

$$\mathbb{E}[\|\nabla_{\theta} J(\pi_{\theta_T})\|] \le \varepsilon + \sqrt{3}(\varepsilon_{reg} + \varepsilon_{data} + \varepsilon_{func})$$

if the hyper-parameters in Alg. 1 and 3 satisfy the following constraints:

$$\begin{split} T = & [\max\{96, \frac{16L_J}{\varepsilon^2}, \frac{16}{(1-\gamma)\varepsilon^2} \sqrt{120\left(2C_{\zeta,\mu}C_{w,Q} + H^2C_Q^2C_W^2\right)}, \frac{864C_{w,Q}d^2}{\varepsilon^2}\}] = O(\varepsilon^{-2});\\ |B| = & [\max\{1, \frac{12\sigma^2}{\varepsilon^2}\}]; \quad |N| = [\frac{96(L^2 + 20C_{w,Q})\sigma^2}{\min\{\frac{\mu_{\zeta}\eta_{\zeta}}{4}, \frac{\mu_{\xi}\eta_{\xi}}{4}\}\varepsilon^2}(\frac{\eta_{\zeta}}{\mu_{\zeta}} + \frac{\eta_{\xi}}{\mu_{\xi}})]; \quad K = c_{oracle}\log(\frac{1}{\beta});\\ \alpha = & 0.9; \quad \beta = \min\{\frac{\varepsilon^2}{L^2}, \frac{(1-\gamma)^2\varepsilon^4}{C_{\zeta,\mu}L^2}, \frac{\alpha}{2}(1-\alpha)^2\}; \quad B_0 = [\frac{8\sigma^2}{\varepsilon^2}]\\ \eta_{\theta} = \min\{\frac{1}{2L_J}, \left(\left[\frac{C_{\zeta,\mu}L^2\beta}{6(1-\beta)} + 120\left(2C_{\zeta,\mu}C_{w,Q} + H^2C_W^2C_Q^2\right)\right]\right)^{-1/2}\} \end{split}$$

where $[\cdot]$ is the upper rounding function, $C_{w,Q} = G^2 L_w^2 C_Q^2 + G^2 C_W^2 L_Q^2$, $C_{\zeta,\mu} = \kappa_{\mu}^2 (\kappa_{\xi} + 1)^2 + \kappa_{\xi}^2 (\kappa_{\mu} + 1)^2$, L_J is defined in Prop. 4.2, η_{ζ} and η_{ξ} satisfy the constraints in Theorem E.1 and c_{oracle} is an independent constant.

Besides, the total gradient computation to obtain θ_{out} should be $|B_0| + |B| \cdot T + |N| \cdot K \cdot T = O(\varepsilon^{-4})$.

We defer the proofs to Appendix D.

5 Conclusion

In this paper, we study two natural optimization strategies for density-ratio based off-policy policy gradients, establish their convergence rates, and characterize the quality of the results. In the future, it will be interesting to extend the results to other settings with milder assumptions, or improve the dependence on ε^{-1} on the convergence rate of our second strategy.

References

- [1] Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. Optimality and approximation with policy gradient methods in markov decision processes. *arXiv preprint arXiv:1908.00261*, 2019.
- [2] Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon: Infinite-horizon off-policy estimation. In Advances in Neural Information Processing Systems, pages 5361–5371, 2018.
- [3] Ruiyi Zhang, Bo Dai, Lihong Li, and Dale Schuurmans. Gendice: Generalized offline estimation of stationary values. In *International Conference on Learning Representations*, 2019.
- [4] Masatoshi Uehara, Jiawei Huang, and Nan Jiang. Minimax weight and q-function learning for off-policy evaluation. arXiv preprint arXiv:1910.12809, 2019.
- [5] Mengjiao Yang, Ofir Nachum, Bo Dai, Lihong Li, and Dale Schuurmans. Off-policy evaluation via the regularized lagrangian. *arXiv preprint arXiv:2007.03438*, 2020.
- [6] Ofir Nachum, Bo Dai, Ilya Kostrikov, Yinlam Chow, Lihong Li, and Dale Schuurmans. Algaedice: Policy gradient from arbitrary experience. arXiv preprint arXiv:1912.02074, 2019.
- [7] Nan Jiang and Jiawei Huang. Minimax confidence interval for off-policy evaluation and policy optimization. *arXiv preprint arXiv:2002.02081*, 2020.
- [8] Luo Luo, Ye Haishan, and Zhang Tong. Stochastic recursive gradient descent ascent for stochastic nonconvex-strongly-concave minimax problems. 2020.
- [9] Ofir Nachum, Yinlam Chow, Bo Dai, and Lihong Li. Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections. In Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, 8-14 December 2019, Vancouver, BC, Canada, pages 2315–2325, 2019.
- [10] Yao Liu, Adith Swaminathan, Alekh Agarwal, and Emma Brunskill. Off-policy policy gradient with state distribution correction. *CoRR*, abs/1904.08473, 2019. URL http://arxiv.org/ abs/1904.08473.
- [11] Thomas Degris, Martha White, and Richard S. Sutton. Off-policy actor-critic. CoRR, abs/1205.4839, 2012. URL http://arxiv.org/abs/1205.4839.
- [12] Ehsan Imani, Eric Graves, and Martha White. An off-policy policy gradient theorem using emphatic weightings. In Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada, pages 96–106, 2018.
- [13] Shangtong Zhang, Bo Liu, Hengshuai Yao, and Shimon Whiteson. Provably convergent two-timescale off-policy actor-critic with function approximation, 2019.
- [14] Matteo Papini, Damiano Binaghi, Giuseppe Canonaco, Matteo Pirotta, and Marcello Restelli. Stochastic variance-reduced policy gradient. *arXiv preprint arXiv:1806.05618*, 2018.
- [15] Pan Xu, Felicia Gao, and Quanquan Gu. An improved convergence analysis of stochastic variance-reduced policy gradient. *arXiv preprint arXiv:1905.12615*, 2019.
- [16] Pan Xu, Felicia Gao, and Quanquan Gu. Sample efficient policy gradient methods with recursive variance reduction. *arXiv preprint arXiv:1909.08610*, 2019.
- [17] Huizhuo Yuan, Xiangru Lian, Ji Liu, and Yuren Zhou. Stochastic recursive momentum for policy gradient methods. arXiv preprint arXiv:2003.04302, 2020.
- [18] F. Huang, Shangqian Gao, Jian Pei, and H. Huang. Momentum-based policy gradient methods. *ArXiv*, abs/2007.06680, 2020.

- [19] Richard S Sutton, David A. McAllester, Satinder P. Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In S. A. Solla, T. K. Leen, and K. Müller, editors, *Advances in Neural Information Processing Systems 12*, pages 1057–1063. MIT Press, 2000.
- [20] Tianyi Lin, Chi Jin, and Michael I Jordan. On gradient descent ascent for nonconvex-concave minimax problems. *arXiv preprint arXiv:1906.00331*, 2019.
- [21] Tianyi Lin, Chi Jin, Michael Jordan, et al. Near-optimal algorithms for minimax optimization. *arXiv preprint arXiv:2002.02417*, 2020.
- [22] Tatjana Chavdarova, Gauthier Gidel, François Fleuret, and Simon Lacoste-Julien. Reducing noise in gan training with variance reduced extragradient. In *Advances in Neural Information Processing Systems*, pages 393–403, 2019.
- [23] Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement learning. In *Proceedings of the 36th International Conference on Machine Learning*, pages 1042–1051, 2019.
- [24] Alessandro Lazaric, Mohammad Ghavamzadeh, and Rémi Munos. Finite-sample analysis of LSTD. In Johannes Fürnkranz and Thorsten Joachims, editors, *Proceedings of the 27th International Conference on Machine Learning (ICML-10), June 21-24, 2010, Haifa, Israel*, pages 615–622. Omnipress, 2010. URL https://icml.cc/Conferences/2010/papers/ 598.pdf.
- [25] Alessandro Lazaric, Mohammad Ghavamzadeh, and Rémi Munos. Finite-sample analysis of least-squares policy iteration. J. Mach. Learn. Res., 13:3041–3074, 2012. URL http: //dl.acm.org/citation.cfm?id=2503339.

A Useful Lemma

Lemma A.1 (Lemma B.2 in [21]). Define

$$\Phi_{\theta}(\zeta) = \min_{\xi \in \Xi} \mathcal{L}^{D}(\theta, \zeta, \xi) \qquad \phi_{\theta}(\zeta) = \arg_{\xi \in \Xi} \min \mathcal{L}^{D}(\theta, \zeta, \xi), \quad for \ \zeta \in \mathbb{R}^{dim(Z)}$$
$$\Psi_{\theta}(\xi) = \max_{\zeta \in Z} \mathcal{L}^{D}(\theta, \zeta, \xi) \qquad \psi_{\theta}(\xi) = \arg_{\zeta \in Z} \max \mathcal{L}^{D}(\theta, \zeta, \xi), \quad for \ \xi \in \mathbb{R}^{dim(\Xi)}$$

Under Assumption A and C, for fixed θ , we have:

(1) The function $\phi_{\theta}(\cdot)$ is $\kappa_{\xi} = \frac{L}{\mu_{\xi}}$ -Lipschitz.

(2) The function $\Phi_{\theta}(\cdot)$ is $2\kappa_{\xi}L = 2\frac{L^2}{\mu_{\xi}}$ -smooth and μ_{ζ} -strongly concave with $\nabla \Phi_{\theta}(\cdot) := \nabla_{\zeta} \mathcal{L}^D(\theta, \zeta, \phi_{\theta}(\zeta)).$

(3) The function $\psi_{\theta}(\cdot)$ is $\kappa_{\zeta} = \frac{L}{\mu_{\zeta}}$ -Lipschitz.

(4) The function $\Psi_{\theta}(\cdot)$ is $2\kappa_{\zeta}L = 2\frac{L^2}{\mu_{\zeta}}$ -smooth and μ_{ξ} -strongly convex with $\nabla\Psi_{\theta}(\cdot) := \nabla_{\xi}\mathcal{L}^D(\theta, \psi_{\theta}(\xi), \xi).$

Remark A.2 (For clarification). According to Danskin's Theorem, in $\nabla \Phi_{\theta}(\cdot) := \nabla_{\zeta} \mathcal{L}^{D}(\theta, \zeta, \phi_{\theta}(\zeta))$, when we compute $\nabla_{\zeta} \mathcal{L}^{D}(\theta, \zeta, \phi_{\theta}(\zeta))$, we treat $\phi_{\theta}(\zeta)$ as a constant, instead of a function w.r.t. ζ . Therefore, for arbitrary ζ', ξ' , based on Assumption A, we always have:

$$\|\nabla \Phi_{\theta}(\cdot) - \nabla_{\zeta} \mathcal{L}^{D}(\theta, \zeta', \xi')\| \le L \|\zeta - \zeta'\| + L \|\phi_{\theta}(\zeta) - \xi'\|$$

We have a similar clarification w.r.t. $\nabla_{\xi} \Psi(\xi)$ *.*

Lemma A.3. For α -strongly-convex function f(x) and β -strongly-concave function g(x) w.r.t. $x \in X$, where $X \subseteq \mathbb{R}^n$ is a convex set. Then, we have

$$\|x - x_f^*\| \le \frac{1}{\alpha} \|\nabla_x f(x)\| \tag{8}$$

$$\frac{\alpha}{2} \|x - x_f^*\|^2 \le f(x) - f(x_f^*) \tag{9}$$

$$\|x - x_g^*\| \le \frac{1}{\beta} \|\nabla_x g(x)\|$$
(10)

$$\frac{\beta}{2} \|x - x_f^*\|^2 \le g(x_g^*) - g(x) \tag{11}$$

where x_f^* and x_g^* the minimum and maximum of f(x) and g(x), respectively.

Proof. Since f(x) is α -strongly-convex, we have

$$(\nabla_x f(x) - \nabla_x f(x_f^*))^\top (x - x_f^*) \ge \alpha ||x - x_f^*||^2$$

$$f(x) \ge f(x_f^*) + \nabla_x f(x_f^*)^\top (x - x_f^*) + \frac{\alpha}{2} ||x - x_f^*||^2$$

Since x_f^* is the minimizer of f(x), we know that

$$\nabla_x f(x_f^*)^\top (x - x_f^*) \ge 0$$

Combining all the above inequalities together and we obtain

$$\|x - x_f^*\|^2 \le \frac{1}{\alpha} \nabla_x f(x)^\top (x - x_f^*) \le \frac{1}{\alpha} \|\nabla_x f(x)\| \|x - x_f^*\|$$
$$f(x) \ge f(x_f^*) + \frac{\alpha}{2} \|x - x_f^*\|^2$$

which implies

$$\|x - x_f^*\| \le \frac{1}{\alpha} \|\nabla_x f(x)\|$$
$$\frac{\alpha}{2} \|x - x_f^*\|^2 \le f(x) - f(x_f^*)$$

By applying the above results for -g(x) which is a β -strongly-convex function and we can complete the proof.

Lemma A.4. For positive definite matrix **A**, and arbitrary $\alpha > 0$, we have:

$$(\mathbf{A}^{\top}\mathbf{A})^{-1}\succ \Big((\alpha\mathbf{I}+\mathbf{A})^{\top}(\alpha\mathbf{I}+\mathbf{A})\Big)^{-1}$$

Proof. Suppose for symmetric matrix A and B, we have the relationship $A \succ B \succ 0$. According to the inverse matrix lemma, we have

$$\mathbf{B}^{-1} - \mathbf{A}^{-1} = \mathbf{B}^{-1} - (\mathbf{B} + (\mathbf{A} - \mathbf{B}))^{-1} = (\mathbf{B} + \mathbf{B}(\mathbf{A} - \mathbf{B})^{-1}\mathbf{B})^{-1}$$

Because $\mathbf{A} \succ \mathbf{B} \succ 0$, we have $(\mathbf{B} + \mathbf{B}(\mathbf{A} - \mathbf{B})^{-1}\mathbf{B})^{-1} \succ 0$, therefore $\mathbf{B}^{-1} \succ \mathbf{A}^{-1}$.

Then, we only need to prove

$$(\alpha \mathbf{I} + \mathbf{A})^{\top} (\alpha \mathbf{I} + \mathbf{A}) \succ \mathbf{A}^{\top} \mathbf{A}$$

We have

$$(\alpha \mathbf{I} + \mathbf{A})^{\top} (\alpha \mathbf{I} + \mathbf{A}) = \alpha^2 \mathbf{I} + \alpha (\mathbf{A} + \mathbf{A}^{\top}) + \mathbf{A}^{\top} \mathbf{A}$$

 \square

Combining $\mathbf{A} = \mathbf{A}^{\top} \succ 0$ and $\alpha > 0$, we can finish the proof.

Lemma A.5 (Non-negative Elements). We use $\mathbf{P}_*^{\pi} = (\mathbf{P}^{\pi})^{\top} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|}$ to denote the transpose of the transition kernel. All the elements in $(\mathbf{I} - \gamma \mathbf{P}_*^{\pi})^{-1}$ are non-negative. Moreover, the element indexed by (s_i, a_j) in row and (s_p, a_q) in column equals to the unnormalized discounted state-action occupancy of (s_i, a_j) starting from (s_p, a_q) and executing π .

Proof. For arbitrary initial state-action distribution vector $\mu_0 \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|\times 1}$, $(\mathbf{I} - \gamma \mathbf{P}_*^{\pi})^{-1}\mu_0$ is a vector whose elements are unnormalized state-action occupancy with μ_0 as initial distribution, which is larger or equal to 0. As a result, by choosing standard basis vector as μ_0 , we can finish the proof.

Proposition 2.4. Denote $\varepsilon_{data} := (2\kappa_{\zeta}\kappa_{\xi} + 2\kappa_{\zeta} + 2\kappa_{\zeta} + \sqrt{2}/2)\sqrt{2\varepsilon'_{data}}$, where ε'_{data} is defined in Definition 2.3, $\kappa_{\zeta} = L/\mu_{\zeta}$, $\kappa_{\xi} = L/\mu_{\xi}$. Under Assumption A and C, we have:

$$\|\nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}(\pi_{\theta}, w, Q) - \nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi_{\theta}, w, Q)\| \leq \varepsilon_{data}$$

Proof. For the simplicity of notation, we give the proof for a fixed π .

Denote $(w_{\mu}^{*}, Q_{\mu}^{*})$ parameterized by $(\zeta_{\mu}^{*}, \xi_{\mu}^{*})$ as $\arg \max_{w \in \mathcal{W}} \arg \min_{Q \in \mathcal{Q}} \mathcal{L}(\pi, w, Q)$ and denote (w^{*}, Q^{*}) parameterized by (ζ^{*}, ζ^{*}) as $\arg \max_{w \in \mathcal{W}} \arg \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi, w, Q)$. First, we try to bound $\zeta^{*} - \zeta_{\mu}^{*}$. We use Q_{w} and Q_{w}^{D} (parameterized by ξ_{w} and ξ_{w}^{D}) as the short notes of $\arg \min_{Q \in \mathcal{Q}} \mathcal{L}(\pi, w, Q)$ and $\arg \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi, w, Q)$, respectively. Then,

$$\begin{aligned} &|\mathcal{L}(\pi, w, Q_w) - \mathcal{L}^D(\pi, w, Q_w^D)| \\ &\leq \max\{\mathcal{L}(\pi, w, Q_w^D) - \mathcal{L}(\pi, w, Q_w), \mathcal{L}^D(\pi, w, Q_w) - \mathcal{L}^D(\pi, w, Q_w^D)\} \leq \varepsilon'_{date} \end{aligned}$$

As a result,

$$\mathcal{L}^{D}(\pi, w^{*}, Q^{*}) - \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi, w^{*}_{\mu}, Q)$$

$$\leq \mathcal{L}^{D}(\pi, w^{*}, Q^{*}) - \min_{Q \in \mathcal{Q}} \mathcal{L}(\pi, w^{*}, Q) + \mathcal{L}(\pi, w^{*}_{\mu}, Q^{*}_{\mu}) - \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi, w^{*}_{\mu}, Q)$$

$$\leq 2\varepsilon'_{data}$$

According to Lemma A.1, $\min_{Q \in \mathcal{Q}} \mathcal{L}^D(\pi, w, Q)$ is μ_{ζ} -strongly concave. Therefore,

$$\|\zeta^* - \zeta^*_{\mu}\| \le \frac{2}{\mu_{\zeta}} \sqrt{\mathcal{L}^D(\pi, w^*, Q^*) - \min_{Q \in \mathcal{Q}} \mathcal{L}^D(\pi, w^*_{\mu}, Q)} \le \frac{2}{\mu_{\zeta}} \sqrt{2\varepsilon'_{data}}$$

Next, we bound $\|\xi^* - \xi^*_{\mu}\|$. For arbitrary $\pi \in \Pi$ and $w \in \mathcal{W}$, we have:

 $\mathcal{L}^{D}(\pi, w, Q_w) - \mathcal{L}^{D}(\pi, w, Q_w^D) \leq \mathcal{L}^{D}(\pi, w, Q_w) - \mathcal{L}(\pi, w, Q_w) + \mathcal{L}(\pi, w, Q_w^D) - \mathcal{L}^{D}(\pi, w, Q_w^D) \leq 2\varepsilon'_{data}$ Since L^{D} is μ_{ξ} strongly-convexity, as a result of Lemma A.3,

$$\|\xi_w - \xi_w^D\| \le \frac{2}{\mu_{\xi}} \sqrt{2\varepsilon'_{data}}$$
(12)

Then, we have

$$\begin{split} \|\xi^{*} - \xi^{*}_{\mu}\| &\leq \|\xi^{*} - \arg\min_{\xi \in \Xi} \mathcal{L}^{D}(\pi, w^{*}_{\mu}, Q_{\xi})\| + \|\arg\min_{\xi \in \Xi} \mathcal{L}^{D}(\pi, w^{*}_{\mu}, Q_{\xi}) - \xi^{*}_{\mu}\| \\ &= \|\xi^{*} - \arg\min_{\xi \in \Xi} \mathcal{L}^{D}(\pi, w^{*}_{\mu}, Q_{\xi})\| + \|\arg\min_{\xi \in \Xi} \mathcal{L}^{D}(\pi, w^{*}_{\mu}, Q_{\xi}) - \arg\min_{\xi \in \Xi} \mathcal{L}(\pi, w^{*}_{\mu}, Q_{\xi})\| \\ &\leq \frac{L}{\mu_{\xi}} \|\zeta^{*} - \zeta^{*}_{\mu}\| + \frac{2}{\mu_{\xi}} \sqrt{2\varepsilon'_{data}} \\ &\leq (\frac{2L}{\mu_{\xi}\mu_{\zeta}} + \frac{2}{\mu_{\xi}}) \sqrt{2\varepsilon'_{data}} \end{split}$$

where in the last but two step, we use Lemma A.1-(1).

As a directly application of Assumption A, we have:

$$\begin{aligned} \|\nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}(\pi_{\theta}, w, Q) - \nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi_{\theta}, w, Q) \| \\ = \|\nabla_{\theta} \mathcal{L}(\pi_{\theta}, w_{\mu}^{*}, Q_{\mu}^{*}) - \nabla_{\theta} \mathcal{L}^{D}(\pi_{\theta}, w_{\mu}^{*}, Q_{\mu}^{*})\| + \|\nabla_{\theta} \mathcal{L}^{D}(\pi_{\theta}, w_{\mu}^{*}, Q_{\mu}^{*}) - \nabla_{\theta} \mathcal{L}^{D}(\pi_{\theta}, w^{*}, Q^{*})\| \\ \leq \sqrt{\varepsilon'_{data}} + L \|\zeta^{*} - \zeta^{*}_{\mu}\| + L \|\xi^{*} - \xi^{*}_{\mu}\| \\ \leq (2\kappa_{\zeta}\kappa_{\xi} + 2\kappa_{\zeta} + 2\kappa_{\xi} + \sqrt{2}/2)\sqrt{2\varepsilon'_{data}} \end{aligned}$$

Proposition 2.6. Under Assumption A and C, for arbitrary $\pi \in \Pi$, we have:

$$\mathbb{E}_{d^{\mu}}[|w_{\mu}^{*}(s,a) - w_{\mathcal{L}}^{\pi}(s,a)|^{2}] \leq \varepsilon_{\mathcal{W}} := 4\frac{\lambda_{\max}^{2}}{\lambda_{Q}\lambda_{w}}\varepsilon_{1} + 2\frac{L_{w}^{2}\lambda_{\max}}{\mu_{\zeta}}\varepsilon_{2}$$
$$\mathbb{E}_{d^{\mu}}[|Q_{\mu}^{*}(s,a) - Q_{\mathcal{L}}^{\pi}(s,a)|^{2}] \leq \varepsilon_{\mathcal{Q}} := 8\frac{\lambda_{\max}^{3}}{\lambda_{Q}^{2}\lambda_{w}}\varepsilon_{1} + (2 + 4\frac{L_{w}^{2}\lambda_{\max}^{2}}{\lambda_{Q}\mu_{\zeta}})\varepsilon_{2}$$

where $(w_{\mu}^{*}, Q_{\mu}^{*})$ denotes the saddle point of $\mathcal{L}(\pi, w, Q)$ constrained by $w, Q \in \mathcal{W} \times \mathcal{Q}$, $(w_{\mathcal{L}}^{\pi}, Q_{\mathcal{L}}^{\pi})$ denotes the saddle point of $\mathcal{L}(\pi, w, Q)$ without any constraint on w and Q, $\lambda_{\max} = \max\{\lambda_Q, \lambda_w\}$, L_w is defined in Assumption A, μ_{ζ} is defined in Assumption C.

Proof. In the following, we will frequently consider two loss functions. The first one is $\mathcal{L}(\pi, w, Q)$ defined in Eq.(1), where w and Q are parameterized by ζ and ξ , respectively, and we will write $(w, Q) \in \mathcal{W} \times \mathcal{Q}$. The second one is $\mathcal{F}(\pi, x, y)$ defined by:

$$\mathcal{F}(\pi, x, y) = (1 - \gamma)(\nu_0^{\pi})^{\top} \mathbf{\Lambda}^{-1/2} y + x^{\top} \left(\mathbf{\Lambda}^{1/2} R - (\mathbf{I} - \gamma \mathbf{\Lambda}^{1/2} P^{\pi} \mathbf{\Lambda}^{-1/2}) y \right) + \frac{\lambda_Q}{2} y^{\top} y - \frac{\lambda_w}{2} x^{\top} x^{\top} x^{\top} y + \frac{\lambda_Q}{2} y^{\top} y - \frac{\lambda_W}{2} x^{\top} x^{\top} y + \frac{\lambda_Q}{2} y^{\top} y - \frac{\lambda_W}{2} x^{\top} x^{\top} y + \frac{\lambda_Q}{2} y^{\top} y - \frac{\lambda_W}{2} x^{\top} x^{\top} y + \frac{\lambda_Q}{2} y^{\top} y - \frac{\lambda_W}{2} x^{\top} y + \frac{\lambda_Q}{2} y^{\top} y - \frac{\lambda_W}{2} x^{\top} x^{\top} y + \frac{\lambda_Q}{2} y^{\top} y - \frac{\lambda_W}{2} x^{\top} y + \frac{\lambda_Q}{2} x^{\top} y + \frac{\lambda_$$

where $(x, y) \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} \times \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$. For simplification, in the following, we will use $\max_x \min_y$ as a short note of $\max_{x \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}} \min_{y \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}}$.

As we can see, the difference between $\mathcal{L}(\pi, w, Q)$ and $\mathcal{F}(\pi, x, y)$ is not only that we don't have any constraint on x and y, but also that we absorb one $\Lambda^{1/2}$ into vector x and y. In another word, for arbitrary π, w, Q , we have

$$\mathcal{L}(\pi, w, Q) = \mathcal{F}(\pi, \Lambda^{1/2} w, \Lambda^{1/2} Q).$$

Obviously, $\mathcal{F}(\pi, x, y)$ is λ_w -strongly-concave- λ_Q -strongly-convex and λ_{\max} -smooth w.r.t. $x, y \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$.

In the following, we use $w_{\mathbb{R}}^*$ parameterized by $\zeta_{\mathbb{R}}^*$ to denote $\arg \max_{w \in \mathcal{W}} \min_y \mathcal{F}(\pi, \Lambda^{1/2}w, y)$. According to Lemma A.1, $\min_y \mathcal{F}(\pi, x, y)$ is a $2\frac{\lambda_{\max}^2}{\lambda_Q}$ -smooth and λ_w -strongly-concave function with gradient $\nabla_x \min_y \mathcal{F}(\pi, x, y)$. Since $\nabla_x \mathcal{F}(\pi, \Lambda^{1/2}w_{\mathcal{L}}^{\pi}, \Lambda^{1/2}Q_{\mathcal{L}}^{\pi}) = 0$, we have,

$$\begin{split} &\frac{\lambda_{w}}{2} \| \mathbf{\Lambda}^{1/2} w_{\mathbb{R}}^{*} - \mathbf{\Lambda}^{1/2} w_{\mathcal{L}}^{\pi} \|^{2} \\ \leq & \mathcal{F}(\pi, \mathbf{\Lambda}^{1/2} w_{\mathcal{L}}^{\pi}, \mathbf{\Lambda}^{1/2} Q_{\mathcal{L}}^{\pi}) - \min_{y} \mathcal{F}(\pi, \mathbf{\Lambda}^{1/2} w_{\mathbb{R}}^{*}, y) \quad (\text{Strong concavity of } \min_{y} \mathcal{F}(\pi, x, y)) \\ = & \mathcal{F}(\pi, \mathbf{\Lambda}^{1/2} w_{\mathcal{L}}^{\pi}, \mathbf{\Lambda}^{1/2} Q_{\mathcal{L}}^{\pi}) - \max_{w \in \mathcal{W}} \min_{y} \mathcal{F}(\pi, \mathbf{\Lambda}^{1/2} w, y) \\ \leq & \mathcal{F}(\pi, \mathbf{\Lambda}^{1/2} w_{\mathcal{L}}^{\pi}, \mathbf{\Lambda}^{1/2} Q_{\mathcal{L}}^{\pi}) - \min_{y} \mathcal{F}(\pi, \mathbf{\Lambda}^{1/2} w_{\zeta^{\pi}}, y) \quad (w_{\zeta^{\pi}} \text{ is defined in Def. 2.5}) \\ \leq & \frac{\lambda_{\max}^{2}}{\lambda_{Q}} \| \mathbf{\Lambda}^{1/2} w_{\zeta^{\pi}} - \mathbf{\Lambda}^{1/2} w_{\mathcal{L}}^{\pi} \|^{2} \quad (\text{Smoothness of } \min_{y} \mathcal{F}(\pi, x, y)) \\ = & \frac{\lambda_{\max}^{2}}{\lambda_{Q}} \| w_{\zeta^{\pi}} - w_{\mathcal{L}}^{\pi} \|_{\mathbf{\Lambda}}^{2} = \frac{\lambda_{\max}^{2}}{\lambda_{Q}} \varepsilon_{1} \quad (\text{see definition of } \varepsilon_{1} \text{ in Def.2.5}) \end{split}$$

which implies

.....

$$\|\mathbf{\Lambda}^{1/2} w_{\mathbb{R}}^* - \mathbf{\Lambda}^{1/2} w_{\mathcal{L}}^{\pi}\|^2 \le 2 \frac{\lambda_{\max}^2}{\lambda_Q \lambda_w} \varepsilon_1 \tag{13}$$

Applying Lemma A.1 for $(w, Q) \in \mathcal{W} \times \mathcal{Q}$, we know $\min_{\xi \in \Xi} \mathcal{L}(\pi, w_{\zeta}, Q_{\xi})$ is μ_{ζ} -strongly-concave w.r.t. ζ . Since ζ^* is the minimizer of $\min_{\xi \in \Xi} \mathcal{L}(\pi, w_{\zeta}, Q_{\xi})$ and Z is a convex set, we have

In the last but two inequality, we use the fact that $\mathcal{F}(\pi, \Lambda^{1/2} w_{\mu}^*, \cdot)$ is λ_{\max} -smooth and $\nabla_y \min_y \mathcal{F}(\pi, \Lambda^{1/2} w_{\mu}^*, Q) = 0$; in the last equality, we use the definition of ε_2 in Def. 2.5-(2). Combing (2) in Assumption A, for arbitrary $s, a \in S \times A$, we have:

$$|w_{\mu}^{*}(s,a) - w_{\mathbb{R}}^{*}(s,a)|^{2} \leq L_{w}^{2} ||\zeta^{*} - \zeta_{\mathcal{R}}^{*}||^{2} \leq \frac{L_{w}^{2} \lambda_{\max}}{\mu_{\zeta}} \varepsilon_{2}$$
(14)

Therefore, as a result of Eq.(13) and Eq.(14):

$$\begin{split} \mathbb{E}_{d^{\mu}}[|w_{\mu}^{*}-w_{\mathcal{L}}^{\pi}|^{2}] \leq & 2\mathbb{E}_{d^{\mu}}[|w_{\mathbb{R}}^{*}-w_{\mathcal{L}}^{\pi}|^{2}] + 2\mathbb{E}_{d^{\mu}}[|w_{\mathbb{R}}^{*}-w_{\mu}^{*}|^{2}] \\ = & 2\|\mathbf{\Lambda}^{1/2}w_{\mathbb{R}}^{*}-\mathbf{\Lambda}^{1/2}w_{\mathcal{L}}^{\pi}\|^{2} + 2\mathbb{E}_{d^{\mu}}[|w_{\mathbb{R}}^{*}-w_{\mu}^{*}|^{2}] \\ \leq & 4\frac{\lambda_{\max}^{2}}{\lambda_{Q}\lambda_{w}}\varepsilon_{1} + 2\frac{L_{w}^{2}\lambda_{\max}}{\mu_{\zeta}}\varepsilon_{2} \end{split}$$

According to Lemma A.1 again, $\arg \min_y \mathcal{F}(\pi, x, y)$ is $\frac{\lambda_{\max}}{\lambda_Q}$ -Lipschitz w.r.t. x, we have

$$\mathbb{E}_{d^{\mu}}[|Q_{\mu}^{*} - Q_{\mathcal{L}}^{\pi}|^{2}] = \|\mathbf{\Lambda}^{1/2}Q_{\mu}^{*} - \mathbf{\Lambda}^{1/2}Q_{\mathcal{L}}^{\pi}\|^{2}$$

$$\leq 2 \underbrace{\|\mathbf{\Lambda}^{1/2}Q_{\mu}^{*} - \arg\min_{y} \mathcal{F}(\pi, \mathbf{\Lambda}^{1/2}w_{\mu}^{*}, Q)\|^{2}}_{bounded \ by \ \varepsilon_{2}} + 2 \|\arg\min_{y} \mathcal{F}(\pi, \mathbf{\Lambda}^{1/2}w_{\mu}^{*}, y) - \mathbf{\Lambda}^{1/2}Q_{\mathcal{L}}^{\pi}\|^{2}$$

$$\leq 2\varepsilon_{2} + 2 \frac{\lambda_{\max}}{\lambda_{Q}} \|\mathbf{\Lambda}^{1/2}w_{\mu}^{*} - \mathbf{\Lambda}^{1/2}w_{\mathcal{L}}^{\pi}\|^{2}$$

$$\leq 8 \frac{\lambda_{\max}^{3}}{\lambda_{Q}^{2}\lambda_{w}} \varepsilon_{1} + (2 + 4 \frac{L_{w}^{2}\lambda_{\max}^{2}}{\lambda_{Q}\mu_{\zeta}})\varepsilon_{2}$$

As a result,

$$\varepsilon_{\mathcal{W}} = 4 \frac{\lambda_{\max}^2}{\lambda_Q \lambda_w} \varepsilon_1 + 2 \frac{L_w^2 \lambda_{\max}}{\mu_\zeta} \varepsilon_2; \quad \varepsilon_{\mathcal{Q}} = 8 \frac{\lambda_{\max}^3}{\lambda_Q^2 \lambda_w} \varepsilon_1 + (2 + 4 \frac{L_w^2 \lambda_{\max}^2}{\lambda_Q \mu_\zeta}) \varepsilon_2$$

B The analysis of Bias

Theorem B.1 (Bias resulting from regularization). Let's rewrite Eq.(1) in a vector-matrix form:

$$\max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}(\pi, w, Q) := (1 - \gamma) (\nu_0^{\pi})^{\top} Q + w^{\top} \mathbf{\Lambda} \Big(R - (I - \gamma \mathbf{P}^{\pi}) Q \Big) + \frac{\lambda_Q}{2} Q^{\top} \mathbf{\Lambda} Q - \frac{\lambda_w}{2} w^{\top} \mathbf{\Lambda} w$$

where ν_0^{π} and \mathbf{P}^{π} denotes the initial state-action distribution and the transition matrix w.r.t. policy π , respectively; $\mathbf{\Lambda} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|}$ denotes the diagonal matrix whose diagonal elements are $d^{\mu}(\cdot, \cdot)$. Denote $(w_{\mathcal{L}}^{\pi}, Q_{\mathcal{L}}^{\pi})$ as the saddle point of $\mathcal{L}(\pi, w, Q)$ without any constraint on w and Q (i.e. $w, Q \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$), then we have:

$$w_{\mathcal{L}}^{\pi} = w^{\pi} + \left(\lambda_{w}\lambda_{Q}I + (I - \gamma\mathbf{P}^{\pi})\mathbf{\Lambda}^{-1}(I - \gamma\mathbf{P}^{\pi}_{*})\mathbf{\Lambda}\right)^{-1}\left(\lambda_{Q}R - \lambda_{Q}\lambda_{w}w^{\pi}\right)$$
$$Q_{\mathcal{L}}^{\pi} = Q^{\pi} - \left(\lambda_{w}\lambda_{Q}I + \mathbf{\Lambda}^{-1}(I - \gamma\mathbf{P}^{\pi}_{*})\mathbf{\Lambda}(I - \gamma\mathbf{P}^{\pi})\right)^{-1}\left(\lambda_{w}\lambda_{Q}Q^{\pi} + \lambda_{w}(1 - \gamma)\mathbf{\Lambda}^{-1}\nu_{0}^{\pi}\right)$$

where $w^{\pi} = \frac{d^{\pi}}{d^{\mu}}$ is the density ratio and Q^{π} is the Q function of π . we use $\mathbf{P}_{*}^{\pi} = (\mathbf{P}^{\pi})^{\top}$ to denote the transpose of the transition matrix.

Proof. Recall the loss function

$$\mathcal{L}(\pi, w, Q) = (1 - \gamma)(\nu_0^{\pi})^{\top} Q + w^{\top} \mathbf{\Lambda} R - w^{\top} \mathbf{\Lambda} (I - \gamma \mathbf{P}^{\pi}) Q + \frac{\lambda_Q}{2} Q^{\top} \mathbf{\Lambda} Q - \frac{\lambda_w}{2} w^{\top} \mathbf{\Lambda} w$$

By taking the derivatives w.r.t. Q, since Λ is invertible, the optimal choice of Q should be:

$$Q = \frac{1}{\lambda_Q} \mathbf{\Lambda}^{-1} ((I - \gamma \mathbf{P}_*^{\pi}) \mathbf{\Lambda} w - (1 - \gamma) \nu_0^{\pi})$$

Plug this result in, and we have

$$\mathcal{L}(\pi, w, Q) = -\frac{1}{2\lambda_Q} \Big((1 - \gamma)\nu_0^{\pi} - (I - \gamma \mathbf{P}_*^{\pi})\mathbf{\Lambda}w \Big)^{\top} \mathbf{\Lambda}^{-1} \Big((1 - \gamma)(\nu_0^{\pi}) - (I - \gamma \mathbf{P}_*^{\pi})\mathbf{\Lambda}w \Big) + w^{\top}\mathbf{\Lambda}R - \frac{\lambda_w}{2}w^{\top}\mathbf{\Lambda}w$$

Taking the derivative w.r.t. w, and set it to 0:

$$0 = \frac{1}{\lambda_Q} \mathbf{\Lambda} (I - \gamma \mathbf{P}^{\pi}) \mathbf{\Lambda}^{-1} \Big((1 - \gamma) (\nu_0^{\pi}) - (I - \gamma \mathbf{P}_*^{\pi}) \mathbf{\Lambda} w \Big) + \mathbf{\Lambda} R - \lambda_w \mathbf{\Lambda} w$$

As a result,

$$\begin{split} w_{\mathcal{L}}^{\pi} &= \left(\lambda_{w}I + \frac{1}{\lambda_{Q}}(I - \gamma \mathbf{P}^{\pi})\mathbf{\Lambda}^{-1}(I - \gamma \mathbf{P}_{*}^{\pi})\mathbf{\Lambda}\right)^{-1} \left(\frac{1}{\lambda_{Q}}(I - \gamma \mathbf{P}^{\pi})\mathbf{\Lambda}^{-1}(1 - \gamma)\nu_{0}^{\pi} + R\right) \\ &= \left(\lambda_{w}\lambda_{Q}I + (I - \gamma \mathbf{P}^{\pi})\mathbf{\Lambda}^{-1}(I - \gamma \mathbf{P}_{*}^{\pi})\mathbf{\Lambda}\right)^{-1} \left((I - \gamma \mathbf{P}^{\pi})\mathbf{\Lambda}^{-1}(I - \gamma \mathbf{P}_{*}^{\pi})\mathbf{\Lambda}\mathbf{\Lambda}^{-1}(I - \gamma \mathbf{P}_{*}^{\pi})^{-1}(1 - \gamma)\nu_{0}^{\pi} + \lambda_{Q}R\right) \\ &= w^{\pi} + \left(\lambda_{w}\lambda_{Q}I + (I - \gamma \mathbf{P}^{\pi})\mathbf{\Lambda}^{-1}(I - \gamma \mathbf{P}_{*}^{\pi})\mathbf{\Lambda}\right)^{-1} \left(\lambda_{Q}R - \lambda_{Q}\lambda_{w}w^{\pi}\right) \end{split}$$

and

$$\begin{aligned} Q_{\mathcal{L}}^{\pi} &= \frac{1}{\lambda_Q} \mathbf{\Lambda}^{-1} \Big((I - \gamma \mathbf{P}_*^{\pi}) \mathbf{\Lambda} w_{\mathcal{L}}^{\pi} - (1 - \gamma) \nu_0^{\pi} \Big) \\ &= \frac{1}{\lambda_Q} \mathbf{\Lambda}^{-1} \Big((I - \gamma \mathbf{P}_*^{\pi}) \mathbf{\Lambda} w_{\mathcal{L}}^{\pi} - (I - \gamma \mathbf{P}_*^{\pi}) \mathbf{\Lambda} w^{\pi} \Big) \\ &= \frac{1}{\lambda_Q} \mathbf{\Lambda}^{-1} (I - \gamma \mathbf{P}_*^{\pi}) \mathbf{\Lambda} \Big(\lambda_Q \lambda_w \mathbf{\Lambda} + \mathbf{\Lambda} (I - \gamma \mathbf{P}^{\pi}) \mathbf{\Lambda}^{-1} (I - \gamma \mathbf{P}_*^{\pi}) \mathbf{\Lambda} \Big)^{-1} \Big(\lambda_Q \mathbf{\Lambda} R - \lambda_Q \lambda_w \mathbf{\Lambda} w^{\pi} \Big) \\ &= \Big(\lambda_w \lambda_Q (I - \gamma \mathbf{P}_*^{\pi})^{-1} \mathbf{\Lambda} + \mathbf{\Lambda} (I - \gamma \mathbf{P}^{\pi}) \Big)^{-1} \Big(\mathbf{\Lambda} R - \lambda_w \mathbf{\Lambda} w^{\pi} \Big) \\ &= \Big(\lambda_w \lambda_Q (I - \gamma \mathbf{P}_*^{\pi})^{-1} \mathbf{\Lambda} + \mathbf{\Lambda} (I - \gamma \mathbf{P}^{\pi}) \Big)^{-1} \Big(\mathbf{\Lambda} (I - \gamma \mathbf{P}^{\pi}) Q^{\pi} - \lambda_w \mathbf{\Lambda} w^{\pi} \Big) \\ &= Q^{\pi} - \Big(\lambda_w \lambda_Q (I - \gamma \mathbf{P}_*^{\pi})^{-1} \mathbf{\Lambda} + \mathbf{\Lambda} (I - \gamma \mathbf{P}^{\pi}) \Big)^{-1} \Big(\lambda_w \lambda_Q (I - \gamma \mathbf{P}_*^{\pi})^{-1} \mathbf{\Lambda} Q^{\pi} + \lambda_w \mathbf{\Lambda} w^{\pi} \Big) \\ &= Q^{\pi} - \Big(\lambda_w \lambda_Q I + \mathbf{\Lambda}^{-1} (I - \gamma \mathbf{P}_*^{\pi}) \mathbf{\Lambda} (I - \gamma \mathbf{P}^{\pi}) \Big)^{-1} \Big(\lambda_w \lambda_Q Q^{\pi} + \lambda_w (1 - \gamma) \mathbf{\Lambda}^{-1} \nu_0^{\pi}) \Big) \end{aligned}$$

Lemma B.2. Under Assumption B:

$$\begin{aligned} \|w^{\pi} - w_{\mathcal{L}}^{\pi}\|_{\mathbf{\Lambda}}^{2} &\leq \frac{C^{2}(\lambda_{Q} + \lambda_{Q}\lambda_{w}C)^{2}}{(1-\gamma)^{4}} \\ \|Q^{\pi} - Q_{\mathcal{L}}^{\pi}\|_{\mathbf{\Lambda}}^{2} &\leq \frac{C^{2}}{(1-\gamma)^{2}}(\frac{\lambda_{w}\lambda_{Q}}{1-\gamma} + \lambda_{w})^{2} \end{aligned}$$

where (w^{π}, Q^{π}) and $(w_{\mathcal{L}}^{\pi}, Q_{\mathcal{L}}^{\pi})$ are defined in Theorem B.1. $||x||_{\mathbf{\Lambda}} = x^{\top} \mathbf{\Lambda} x$ denotes the norm of column vector x weighted by $\mathbf{\Lambda}$.

Proof. From Theorem B.1, we have

$$w_{\mathcal{L}}^{\pi} = w^{\pi} + \left(\lambda_{w}\lambda_{Q}I + (I - \gamma\mathbf{P}^{\pi})\mathbf{\Lambda}^{-1}(I - \gamma\mathbf{P}^{\pi}_{*})\mathbf{\Lambda}\right)^{-1} \left(\lambda_{Q}R - \lambda_{Q}\lambda_{w}w^{\pi}\right)$$
$$Q_{\mathcal{L}}^{\pi} = Q^{\pi} - \left(\lambda_{w}\lambda_{Q}I + \mathbf{\Lambda}^{-1}(I - \gamma\mathbf{P}^{\pi}_{*})\mathbf{\Lambda}(I - \gamma\mathbf{P}^{\pi})\right)^{-1} \left(\lambda_{w}\lambda_{Q}Q^{\pi} + \lambda_{w}(1 - \gamma)\mathbf{\Lambda}^{-1}\nu_{0}^{\pi}\right)$$

We use $\mathbf{1} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times 1}$ to denote a vector whose all elements are 1. Then, we have

$$\begin{split} \|w^{\pi} - w_{\mathcal{L}}^{\pi}\|_{\mathbf{\Lambda}}^{2} &= \|\left(\lambda_{w}\lambda_{Q}I + (I - \gamma\mathbf{P}^{\pi})\mathbf{\Lambda}^{-1}(I - \gamma\mathbf{P}^{\pi}_{*})\mathbf{\Lambda}\right)^{-1}\left(\lambda_{Q}R - \lambda_{Q}\lambda_{w}w^{\pi}\right)\|_{\mathbf{\Lambda}}^{2} \\ &= \|\left(\lambda_{w}\lambda_{Q}I + \mathbf{\Lambda}^{1/2}(I - \gamma\mathbf{P}^{\pi})\mathbf{\Lambda}^{-1}(I - \gamma\mathbf{P}^{\pi}_{*})\mathbf{\Lambda}^{1/2}\right)^{-1}\mathbf{\Lambda}^{1/2}\left(\lambda_{Q}R - \lambda_{Q}\lambda_{w}w^{\pi}\right)\|^{2} \\ &\leq \|\mathbf{\Lambda}^{-1/2}(I - \gamma\mathbf{P}^{\pi}_{*})^{-1}\mathbf{\Lambda}(I - \gamma\mathbf{P}^{\pi})^{-1}\left(\lambda_{Q}R - \lambda_{Q}\lambda_{w}w^{\pi}\right)\|^{2} \\ &= \|\mathbf{\Lambda}^{-1/2}(I - \gamma\mathbf{P}^{\pi}_{*})^{-1}\mathbf{\Lambda}\widetilde{Q}^{\pi}\|^{2} \\ &\leq \frac{(\lambda_{Q} + \lambda_{Q}\lambda_{w}C)^{2}}{(1 - \gamma)^{2}}\|\mathbf{\Lambda}^{-1}(I - \gamma\mathbf{P}^{\pi}_{*})^{-1}\mathbf{\Lambda}\mathbf{1}\|_{\mathbf{\Lambda}}^{2} \\ &= \frac{(\lambda_{Q} + \lambda_{Q}\lambda_{w}C)^{2}}{(1 - \gamma)^{2}}\|\mathbf{\Lambda}^{-1}(I - \gamma\mathbf{P}^{\pi}_{*})^{-1}d^{\mu}\|_{\mathbf{\Lambda}}^{2} \\ &= \frac{(\lambda_{Q} + \lambda_{Q}\lambda_{w}C)^{2}}{(1 - \gamma)^{4}}\|w_{d^{\mu}}^{\pi}\|_{\mathbf{\Lambda}}^{2} \leq \frac{C^{2}(\lambda_{Q} + \lambda_{Q}\lambda_{w}C)^{2}}{(1 - \gamma)^{4}} \end{split}$$

where in the first inequality, we use Lemma A.4; in the third equality, we use \tilde{Q}^{π} to denote the Q function after replacing true rewards with $\lambda_Q R - \lambda_Q \lambda_w w^{\pi}$; in the second inequality, we use Lemma A.5 and the result that $|\lambda_Q R - \lambda_Q \lambda_w w^{\pi}| \leq \lambda_Q + \lambda_Q \lambda_w C$ given Assumption B; in the last inequality, we use Assumption B again. Similarly,

$$\begin{split} \|Q^{\pi} - Q_{\mathcal{L}}^{\pi}\|_{\mathbf{A}}^{2} &\leq \|\left(\lambda_{w}\lambda_{Q}I + \mathbf{\Lambda}^{-1}(I - \gamma\mathbf{P}_{*}^{\pi})\mathbf{\Lambda}(I - \gamma\mathbf{P}^{\pi})\right)^{-1}\left(\lambda_{w}\lambda_{Q}Q^{\pi} + \lambda_{w}(1 - \gamma)\mathbf{\Lambda}^{-1}\nu_{0}^{\pi}\right)\right)\|_{\mathbf{A}}^{2} \\ &= \|\left(\lambda_{Q}\lambda_{w}I + \mathbf{\Lambda}^{-1/2}(I - \gamma\mathbf{P}_{*}^{\pi})\mathbf{\Lambda}(I - \gamma\mathbf{P}^{\pi})\mathbf{\Lambda}^{-1/2}\right)^{-1}\mathbf{\Lambda}^{1/2}\left(\lambda_{Q}\lambda_{w}Q^{\pi} + \lambda_{w}(1 - \gamma)\mathbf{\Lambda}^{-1}\nu_{0}^{\pi}\right)\right)\|^{2} \\ &\leq \|\mathbf{\Lambda}^{1/2}(I - \gamma\mathbf{P}^{\pi})^{-1}\mathbf{\Lambda}^{-1}(I - \gamma\mathbf{P}_{*}^{\pi})^{-1}\left(\lambda_{w}\lambda_{Q}\mathbf{\Lambda}Q^{\pi} + \lambda_{w}(1 - \gamma)\nu_{0}^{\pi}\right)\right)\|^{2} \\ &= \|\lambda_{w}\lambda_{Q}\mathbf{\Lambda}^{1/2}(I - \gamma\mathbf{P}^{\pi})^{-1}\mathbf{\Lambda}^{-1}(I - \gamma\mathbf{P}_{*}^{\pi})^{-1}\mathbf{\Lambda}Q^{\pi} + \lambda_{w}\mathbf{\Lambda}^{1/2}(I - \gamma\mathbf{P}^{\pi})^{-1}w^{\pi})\|^{2} \\ &\leq \|\frac{\lambda_{w}\lambda_{Q}}{1 - \gamma}\mathbf{\Lambda}^{1/2}(I - \gamma\mathbf{P}^{\pi})^{-1}\mathbf{\Lambda}^{-1}(I - \gamma\mathbf{P}_{*}^{\pi})^{-1}\mathbf{\Lambda}\mathbf{1} + \lambda_{w}\mathbf{\Lambda}^{1/2}(I - \gamma\mathbf{P}^{\pi})^{-1}w^{\pi})\|^{2} \\ &\leq \|(I - \gamma\mathbf{P}^{\pi})^{-1}(\frac{\lambda_{w}\lambda_{Q}}{1 - \gamma}w_{d^{\mu}}^{\pi} + \lambda_{w}w^{\pi})\|_{\mathbf{\Lambda}}^{2} \\ &\leq \frac{C^{2}}{(1 - \gamma)^{2}}(\frac{\lambda_{w}\lambda_{Q}}{1 - \gamma} + \lambda_{w})^{2} \end{split}$$

where in the last but third inequality, we use Lemma A.5 and the fact that w^{π} is also non-negative.

Lemma B.3. Under Assumption B, for arbitrary function f(s, a),

$$(1-\gamma)\mathbb{E}_{s_0 \sim \nu_0, a_0 \sim \pi}[f(s_0, a_0)] + \gamma \mathbb{E}_{s, a, s' \sim d^{\mu}, a' \sim \pi}[w^{\pi}(s, a)f(s', a')] = \mathbb{E}_{d^{\mu}}[w^{\pi}(s, a)f(s, a)]$$
(15)

$$\gamma \mathbb{E}_{s,a,s' \sim d^{\mu},a' \sim \pi} [f^2(s',a')] \le \frac{1}{1-\gamma} \mathbb{E}_{s,a \sim d^{\pi}_{d^{\mu}}} [f^2(s,a)] \le \frac{C}{1-\gamma} \mathbb{E}_{s,a \sim d^{\mu}} [f^2(s,a)]$$
(16)

where $d_{d^{\mu}}^{\pi} := (1 - \gamma) \mathbb{E}_{\tau \sim \pi, s_0, a_0 \sim d^{\mu}(\cdot, \cdot)} [\sum_{t=0}^{\infty} \gamma^t p(s_t = s, a_t = a)]$ is the normalized discounted state-action occupancy by treating $d^{\mu}(\cdot, \cdot)$ as initial distribution; $s, a, s' \sim d^{\mu}, a' \sim \pi$ is a short note of $s, a \sim d^{\mu}, s' \sim P(s'|s, a), a' \sim \pi(\cdot|s')$.

Proof. Eq.(15) can be proved by the equation:

$$d^{\pi}(s,a) = (1-\gamma)\nu_0(s)\pi(a|s) + \gamma \sum_{s',a'} p(s|s',a')d^{\pi}(s',a')\pi(a|s)$$

For Eq.(16), the first step is because $\gamma \sum_{s',a'} d^{\mu}(s',a')p(s|s',a')\pi(a|s) \leq \frac{1}{1-\gamma}d^{\pi}_{d^{\mu}}(s,a)$, and the second step is the result of Assumption B.

Theorem 2.7. [*Bias*] Under Assumption A, B, C, given arbitrary $\theta \in \Theta$, we have

$$\|\nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi_{\theta}, w, Q) - \nabla_{\theta} J(\pi_{\theta})\| \le \varepsilon_{reg} + \varepsilon_{func} + \varepsilon_{data}$$

where ε_{data} is defined in Proposition 2.4, and

$$\varepsilon_{func} = \frac{G}{1-\gamma} \Big(\sqrt{C\varepsilon_{\mathcal{Q}}} + C_{\mathcal{W}} \sqrt{\frac{\gamma\varepsilon_{\mathcal{Q}}C}{1-\gamma}} + \sqrt{\frac{\gamma\varepsilon_{\mathcal{Q}}\varepsilon_{\mathcal{W}}C}{1-\gamma}} + \gamma C_{\mathcal{Q}} \sqrt{\varepsilon_{\mathcal{W}}} \Big)$$
(5.1)

(ε_W and ε_Q defined in Prop. 2.6)

$$\varepsilon_{reg} = \frac{G}{1-\gamma} \Big(\frac{C^2}{(1-\gamma)} \big(\frac{\lambda_w \lambda_Q}{1-\gamma} + \lambda_w \big) + \frac{\gamma C (\lambda_Q + \lambda_Q \lambda_w C)}{(1-\gamma)^3} + \frac{C^2 (\lambda_Q + \lambda_Q \lambda_w C)}{(1-\gamma)^3} \big(\frac{\lambda_w \lambda_Q}{1-\gamma} + \lambda_w \big) \sqrt{\frac{\gamma C}{1-\gamma}} \Big)$$

Proof. Firstly, by applying the triangle inequality:

t

$$\begin{aligned} \|\nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi_{\theta}, w, Q) - \nabla_{\theta} J(\pi_{\theta})\| &\leq \underbrace{\|\nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}^{D}(\pi_{\theta}, w, Q) - \nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}(\pi_{\theta}, w, Q)\|}_{Bounded in Assumption 2.3} \\ &+ \underbrace{\|\nabla_{\theta} \max_{w} \min_{Q} \mathcal{L}(\pi_{\theta}, w, Q) - \nabla_{\theta} \max_{w \in \mathcal{W}} \min_{Q \in \mathcal{Q}} \mathcal{L}(\pi_{\theta}, w, Q)\|}_{t_{1}} \\ &+ \underbrace{\|\nabla_{\theta} J(\pi_{\theta}) - \nabla_{\theta} \max_{w} \min_{Q} \mathcal{L}(\pi_{\theta}, w, Q)\|}_{t_{2}} \end{aligned}$$

where we use $\max_{w} \min_{Q}$ as a short note of $\max_{w \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}} \min_{Q \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}}$.

In the following, we again use $(w_{\mathcal{L}}^{\pi_{\theta}}, Q_{\mathcal{L}}^{\pi_{\theta}})$ to denote the saddle point of $\mathcal{L}(\pi_{\theta}, w, Q)$ without any constraint on w and Q, and use $(w_{\mu}^{*}, Q_{\mu}^{*})$ to denote the saddle point of $\mathcal{L}(\pi_{\theta}, w, Q)$. Next, we upper bound t_{1} and t_{2} one by one. For simplicity, we use $s, a, s' \sim d^{\mu}, a' \sim \pi_{\theta}$ as a short note of $s, a \sim d^{\mu}, s' \sim P(s'|s, a), a' \sim \pi_{\theta}(\cdot|s')$.

Upper bound t_1 With misspecification Definition 2.5, we can easily bound t_1 :

$$\begin{split} & 1 = \|\nabla_{\theta}\mathcal{L}(\pi_{\theta}, w_{\mu}^{*}, Q_{\mu}^{*}) - \nabla_{\theta}\mathcal{L}(\pi_{\theta}, w_{\mathcal{L}^{\theta}}^{\pi_{\theta}}, Q_{\mathcal{L}^{\theta}}^{\pi_{\theta}})\| \\ & \leq \frac{1}{1 - \gamma} \|(1 - \gamma)\mathbb{E}_{\nu_{0}^{\pi_{\theta}}}[\left(Q_{\mu}^{*}(s_{0}, a_{0}) - Q_{\mathcal{L}}^{\pi_{\theta}}(s_{0}, a_{0})\right)\nabla_{\theta}\log \pi_{\theta}(a_{0}|s_{0})]\| \\ & + \frac{\gamma}{1 - \gamma} \|\mathbb{E}_{s,a,s' \sim d^{\mu},a' \sim \pi}[w_{\mu}^{*}(s, a)\left(Q_{\mu}^{*}(s', a') - Q_{\mathcal{L}}^{\pi_{\theta}}(s', a')\right)\nabla_{\theta}\log \pi(a'|s')]\| \\ & + \frac{\gamma}{1 - \gamma} \|\mathbb{E}_{s,a,s' \sim d^{\mu},a' \sim \pi}[(w_{\mu}^{*}(s, a) - w_{\mathcal{L}}^{\pi_{\theta}}(s, a))\left(Q_{\mu}^{*}(s', a') - Q_{\mathcal{L}}^{\pi_{\theta}}(s', a')\right)\nabla_{\theta}\log \pi(a'|s')]\| \\ & + \frac{\gamma}{1 - \gamma} \|\mathbb{E}_{s,a,s' \sim d^{\mu},a' \sim \pi_{\theta}}[(w_{\mu}^{*}(s, a) - w_{\mathcal{L}}^{\pi_{\theta}}(s, a))Q_{\mu}^{*}(s', a') - Q_{\mathcal{L}}^{\pi_{\theta}}(s', a')\right)\nabla_{\theta}\log \pi(a'|s')]\| \\ & \leq \frac{G}{1 - \gamma}\mathbb{E}_{\nu_{0}^{\pi_{\theta}}}[|Q_{\mu}^{*}(s, a) - Q_{\mathcal{L}}^{\pi_{\theta}}(s, a)|] + \frac{\gamma C_{W}G}{1 - \gamma}\mathbb{E}_{s,a,s' \sim d^{\mu},a' \sim \pi_{\theta}}[|Q_{\mu}^{*}(s', a') - Q_{\mathcal{L}}^{\pi_{\theta}}(s', a')|] \\ & \quad ((1 - \gamma)\nu_{0}^{\pi}(s, a) \leq d^{\pi}(s, a) \leq Cd^{\mu}(s, a)) \\ & \quad + \frac{\gamma G}{1 - \gamma}\mathbb{E}_{s,a,s' \sim d^{\mu},a' \sim \pi_{\theta}}[|(w_{\mu}^{*}(s, a) - w_{\mathcal{L}}^{\pi_{\theta}}(s, a))]\left(Q_{\mu}^{*}(s', a') - Q_{\mathcal{L}}^{\pi_{\theta}}(s', a')\right)|] \\ & \quad + \frac{\gamma C_{Q}Q}{1 - \gamma}\mathbb{E}_{s,a,s' \sim d^{\mu},a' \sim \pi_{\theta}}[|(w_{\mu}^{*}(s, a) - w_{\mathcal{L}}^{\pi_{\theta}}(s, a))|] \\ & \quad \leq \frac{G}{1 - \gamma}\sqrt{\mathbb{E}_{\nu_{0}^{\pi_{\theta}}}[|Q_{\mu}^{*}(s, a) - Q_{\mathcal{L}}^{\pi_{\theta}}(s, a)|^{2}]} + \frac{\gamma C_{W}G}{1 - \gamma}\sqrt{\mathbb{E}_{s,a,s' \sim d^{\mu},a' \sim \pi_{\theta}}[|Q_{\mu}^{*}(s', a') - Q_{\mathcal{L}}^{\pi_{\theta}}(s', a')|^{2}]} \\ & \quad + \frac{\gamma G}{1 - \gamma}\sqrt{\mathbb{E}_{d^{\mu}}[|w_{\mathcal{L}}^{\pi_{\theta}}(s, a) - w_{\mathcal{L}}^{\pi_{\theta}}(s, a)|^{2}]} + \frac{\gamma C_{W}G}{1 - \gamma}\sqrt{\mathbb{E}_{d^{\mu}}[|w_{\mathcal{L}}^{\pi_{\theta}}(s, a) - w_{\mathcal{L}}^{\pi_{\theta}}(s, a)|^{2}]} \\ & \quad + \frac{\gamma C_{Q}Q}{1 - \gamma}\sqrt{\mathbb{E}_{d^{\mu}}[|w_{\mathcal{L}}^{\pi_{\theta}}(s, a) - w_{\mathcal{L}}^{\pi_{\theta}}(s, a)|^{2}]} \\ \end{array}$$

$$\leq \frac{G}{1-\gamma} \sqrt{C\mathbb{E}_{d^{\mu}}[|Q_{\mu}^{*}(s,a) - Q_{\mathcal{L}}^{\pi_{\theta}}(s,a)|^{2}]} + \frac{C_{\mathcal{W}}G}{1-\gamma} \sqrt{\frac{\gamma C}{1-\gamma}} \mathbb{E}_{d^{\mu}}[|Q_{\mu}^{*}(s,a) - Q_{\mathcal{L}}^{\pi_{\theta}}(s,a)|^{2}]} \\ + \frac{G}{1-\gamma} \sqrt{\frac{\gamma C}{1-\gamma}} \mathbb{E}_{d^{\mu}}[|w_{\mathcal{L}}^{\pi_{\theta}}(s,a) - w_{\mu}^{*}(s,a)|^{2}] \mathbb{E}_{d^{\mu}}[|Q^{\pi_{\theta}}(s,a) - Q_{\mathcal{L}}^{\pi_{\theta}}(s,a)|^{2}]]} \\ + \frac{\gamma C_{\mathcal{Q}}G}{1-\gamma} \sqrt{\mathbb{E}_{d^{\mu}}[|w_{\mu}^{*}(s,a) - w_{\mathcal{L}}^{\pi_{\theta}}(s,a)|^{2}]} \\ \leq \frac{G}{1-\gamma} \Big(\sqrt{C\varepsilon_{\mathcal{Q}}} + C_{\mathcal{W}} \sqrt{\frac{\gamma\varepsilon_{\mathcal{Q}}C}{1-\gamma}} + \sqrt{\frac{\gamma\varepsilon_{\mathcal{Q}}\varepsilon_{\mathcal{W}}C}{1-\gamma}} + \gamma C_{\mathcal{Q}} \sqrt{\varepsilon_{\mathcal{W}}}\Big)$$

In the last equation, we first use Eq.(16) in Lemma B.3, and then apply Proposition 2.6.

Upper bound t_2 Similarly, we can give a bound for t_2 :

$$\begin{split} t_{2} &= \|\nabla_{\theta} J(\pi_{\theta}) - \nabla_{\theta} \mathcal{L}(\pi_{\theta}, w_{L}^{\pi_{\theta}}, Q_{L}^{\pi_{\theta}}))\| \\ &\leq \frac{1}{1-\gamma} \|(1-\gamma) \mathbb{E}_{v_{0}^{\pi_{\theta}}}[\left(Q^{\pi_{\theta}}(s_{0}, a_{0}) - Q_{L}^{\pi_{\theta}}(s_{0}, a_{0})\right) \nabla_{\theta} \log \pi_{\theta}(a_{0}|s_{0})] \\ &+ \gamma \mathbb{E}_{d^{\mu}}[w^{\pi_{\theta}}(s, a) \left(Q^{\pi_{\theta}}(s', a') - Q_{L}^{\pi_{\theta}}(s', a')\right) \nabla_{\theta} \log \pi(a'|s')]\| \\ &+ \frac{\gamma}{1-\gamma} \|\mathbb{E}_{d^{\mu}}[[(w^{\pi_{\theta}}(s, a) - w_{L}^{\pi_{\theta}}(s, a)) \left(Q^{\pi_{\theta}}(s', a') - Q_{L}^{\pi_{\theta}}(s', a')\right) \nabla_{\theta} \log \pi(a'|s')]\| \\ &+ \frac{\gamma}{1-\gamma} \|\mathbb{E}_{d^{\mu}}[w^{\pi_{\theta}}(s, a) - w_{L}^{\pi_{\theta}}(s, a)] Q^{\pi_{\theta}}(s', a') - Q_{L}^{\pi_{\theta}}(s', a')\right) \nabla_{\theta} \log \pi(a'|s')]\| \\ &= \frac{1}{1-\gamma} \|\mathbb{E}_{d^{\mu}}[w^{\pi_{\theta}}(s, a) \left(Q^{\pi_{\theta}}(s, a) - Q_{L}^{\pi_{\theta}}(s, a)\right) \nabla_{\theta} \log \pi(a|s)]\| \quad (\text{Eq.(15) in Lemma B.3)} \\ &+ \frac{\gamma}{1-\gamma} \|\mathbb{E}_{s,a,s' \sim d^{\mu},a' \sim \pi_{\theta}}[(w^{\pi_{\theta}}(s, a) - w_{L}^{\pi_{\theta}}(s, a)) \left(Q^{\pi_{\theta}}(s', a') - Q_{L}^{\pi_{\theta}}(s', a')\right) \nabla_{\theta} \log \pi(a'|s')]\| \\ &\leq \frac{CG}{1-\gamma} \mathbb{E}_{d^{\mu}}[Q^{\pi_{\theta}}(s, a) - Q_{L}^{\pi_{\theta}}(s, a)] \\ &+ \frac{\gamma G}{1-\gamma} \mathbb{E}_{s,a,s' \sim d^{\mu},a' \sim \pi_{\theta}}[(w^{\pi_{\theta}}(s, a) - w_{L}^{\pi_{\theta}}(s, a)) \left(Q^{\pi_{\theta}}(s', a') - Q_{L}^{\pi_{\theta}}(s', a')\right)|] \\ &+ \frac{\gamma G}{1-\gamma} \sqrt{\mathbb{E}_{d^{\mu}}[|Q^{\pi_{\theta}}(s, a) - Q_{L}^{\pi_{\theta}}(s, a)]} \\ &= \frac{CG}{1-\gamma} \sqrt{\mathbb{E}_{d^{\mu}}[|Q^{\pi_{\theta}}(s, a) - Q_{L}^{\pi_{\theta}}(s, a)]^{2}} + \frac{\gamma G}{(1-\gamma)^{2}} \sqrt{\mathbb{E}_{d^{\mu}}[|(w^{\pi_{\theta}}(s, a) - w_{L}^{\pi_{\theta}}(s, a)]^{2}} \\ &+ \frac{\gamma G}{1-\gamma} \sqrt{\mathbb{E}_{d^{\mu}}[|Q^{\pi_{\theta}}(s, a) - w^{\pi_{\theta}}(s, a)]^{2}} \\ &= \frac{CG}{1-\gamma} \sqrt{\mathbb{E}_{d^{\mu}}[|Q^{\pi_{\theta}}(s, a) - w^{\pi_{\theta}}(s, a)]^{2}} \\ &+ \frac{G}{1-\gamma} \sqrt{\frac{2}{\mathbb{E}_{d^{\mu}}}[|w_{L}^{\pi_{\theta}}(s, a) - w^{\pi_{\theta}}(s, a)]^{2}} \\ &= \frac{CG}{1-\gamma} \sqrt{\frac{2}{\mathbb{E}_{$$

$$\leq \frac{G}{1-\gamma} \Big(\frac{C^2}{(1-\gamma)} \Big(\frac{\lambda_w \lambda_Q}{1-\gamma} + \lambda_w \Big) + \frac{\gamma C(\lambda_Q + \lambda_Q \lambda_w C)}{(1-\gamma)^3} + \frac{C^2(\lambda_Q + \lambda_Q \lambda_w C)}{(1-\gamma)^3} \Big(\frac{\lambda_w \lambda_Q}{1-\gamma} + \lambda_w \Big) \sqrt{\frac{\gamma C}{1-\gamma}} \Big)$$

Г	-	٦	
L			

B.1 Importance of the Regularization

Here we want to highlight that the additional regularization terms on Q and w are crucial. For example, suppose $Q^{\pi} \in Q$ and $w^{\pi} \in W$ for some policy π , if $\lambda_w = \lambda_Q = 0$, we have

$$\begin{aligned} \forall \zeta \in Z, \quad \nabla_{\zeta} \mathcal{L}^{D}(\pi_{\theta}, w_{\zeta}, Q^{\pi}) &= \nabla_{\zeta} (1 - \gamma) \mathbb{E}_{s_{0} \sim \nu_{0}^{D}}[Q^{\pi}(s_{0}, \pi)] = 0 \\ \forall \xi \in \Xi, \quad \nabla_{\xi} \mathcal{L}^{D}(\pi_{\theta}, w^{\pi}, Q_{\xi}) &= \nabla_{\xi} \mathbb{E}_{w^{\pi/\mu}}[r] = 0 \end{aligned}$$

which means $Q = Q^{\pi}$ (or $w = w^{\pi/\mu}$) can result in that the gradient w.r.t. ζ (or ξ) vanishes to 0, and it's impossible that \mathcal{L}^D is a strongly-concave-strongly-convex function.

C Missing Examples and Proofs in Section 3

C.1 Missing proofs

Theorem 3.1. [Equivalence Between Stationary Points] Under Assumption A, C and D, suppose there exists a $\theta \in \Theta$ s.t. $\|\nabla_{\theta} \max_{\zeta \in Z} \min_{\xi \in \Xi} \mathcal{L}^{D}(\theta, \zeta, \xi)\| = 0$ and there is an Algorithm provides us one stationary point $(\theta_T, \zeta_T, \xi_T)$ of the non-concave-strongly-convex problem $\max_{\theta,\zeta} \min_{\xi} \mathcal{L}^{D}(\theta, \zeta, \xi)$ after running T iterations, which statisfying the following conditions in expectation over the randomness of algorithm.

$$\mathbb{E}[\|\nabla_{\theta,\zeta}\mathcal{L}^{D}(\theta_{T},\zeta_{T},\phi_{\theta_{T}}(\zeta_{T}))\|] = \mathbb{E}[\|\nabla_{\theta}\mathcal{L}^{D}(\theta_{T},\zeta_{T},\phi_{\theta_{T}}(\zeta_{T}))\| + \|\nabla_{\zeta}\mathcal{L}^{D}(\theta_{T},\zeta_{T},\phi_{\theta_{T}}(\zeta_{T}))\|] \le \frac{\varepsilon}{(\kappa_{\varepsilon}+1)(\kappa_{\zeta}+1)}$$
(6)

where $\phi_{\theta}(\zeta) = \arg \min_{\xi \in \Xi} \mathcal{L}^{D}(\theta, \zeta, \xi)$. Then, we have

$$\mathbb{E}[\|\nabla_{\theta} J(\pi_{\theta_T})\|] \le \varepsilon + \varepsilon_{data} + \varepsilon_{func} + \varepsilon_{reg}$$

Proof. First of all, as a results of Assumption A, C and D and the condition that $\|\nabla_{\theta} \max_{\zeta \in Z} \min_{\xi \in \Xi} \mathcal{L}^{D}(\theta, \zeta, \xi)\| = 0$ for some $\theta \in \Theta$, we know there must exists $\theta_{T} \in \Theta$ and $\zeta_{T} \in Z$ which can satisfy Eq.(6). Therefore, it's possible for an algorithm to return us a (θ_{T}, ζ_{T}) satisfy Eq.(6).

Next, suppose we already have Eq.(6), it implies that

$$\max\{\mathbb{E}[\|\nabla_{\theta}\mathcal{L}^{D}(\theta_{T},\zeta_{T},\phi_{\theta_{T}}(\zeta_{T}))\|],\mathbb{E}[\|\nabla_{\zeta}\mathcal{L}^{D}(\theta_{T},\zeta_{T},\phi_{\theta_{T}}(\zeta_{T}))\|]\} \le \frac{\varepsilon}{(\kappa_{\xi}+1)(\kappa_{\zeta}+1)}$$
(17)

We can upper bounded $\mathbb{E}[\|\nabla_{\theta} J(\pi_{\theta_T})\|]$ with the triangle inequality:

$$\mathbb{E}[\|\nabla_{\theta}J(\pi_{\theta_{T}})\|] \leq \underbrace{\mathbb{E}[\|\nabla_{\theta}\mathcal{L}^{D}(\theta_{T},\zeta_{T},\phi_{\theta_{T}}(\zeta_{T}))\|]}_{Bounded in Eq.(17)} + \mathbb{E}[\|\nabla_{\theta}\mathcal{L}^{D}(\theta_{T},\zeta^{*},\xi^{*}) - \nabla_{\theta}J(\pi_{\theta_{T}})\|]}_{Bounded in Theorem 2.7} \leq \frac{\varepsilon}{(\kappa_{\xi}+1)(\kappa_{\zeta}+1)} + \varepsilon_{func} + \varepsilon_{reg} + \varepsilon_{data} + \mathbb{E}[\|\nabla_{\theta}\mathcal{L}^{D}(\theta_{T},\zeta^{*},\xi^{*}) - \nabla_{\theta}\mathcal{L}^{D}(\theta_{T},\zeta_{T},\phi_{\theta_{T}}(\zeta_{T})))\|]$$

where we use ζ^*, ξ^* to denote the saddle-point of $\max_{\zeta \in \mathbb{Z}} \min_{\xi \in \Xi} \mathcal{L}^D(\theta_T, \zeta, \xi)$; in the last inequality we use Eq.17 and Theorem 2.7.

Next, we try to bound the last term. According to the definition, ζ^* is also the maximum of function $\Phi_{\theta_T}(\cdot) = \min_{\xi \in \Xi} \mathcal{L}^D(\theta_T, \cdot, \xi)$ defined in Lemma A.1. Applying Property (2) in Lemma A.1, (10) in Lemma A.3 and inequality (17), we obtain that

$$\|\zeta_T - \zeta^*\| \le \frac{1}{\mu_{\zeta}} \|\Phi_{\theta_T}(\zeta_T)\| = \frac{1}{\mu_{\zeta}} \|\nabla_{\zeta} \mathcal{L}^D(\theta_T, \zeta_T, \phi_{\theta_T}(\zeta_T))\| \le \frac{\varepsilon}{\mu_{\zeta}(\kappa_{\xi} + 1)(\kappa_{\zeta} + 1)}$$

Then we can bound:

$$\begin{aligned} \|\nabla_{\theta} \mathcal{L}^{D}(\theta_{T}, \zeta^{*}, \xi^{*}) - \nabla_{\theta} \mathcal{L}^{D}(\theta_{T}, \zeta_{T}, \phi_{\theta_{T}}(\zeta_{T}))\| \\ \leq L \|\zeta_{T} - \zeta^{*}\| + L \|\xi^{*} - \phi_{\theta_{T}}(\zeta_{T}))\| &= L \|\zeta_{T} - \zeta^{*}\| + L \|\phi_{\theta_{T}}(\zeta^{*}) - \phi_{\theta_{T}}(\zeta_{T}))\| \\ \leq (L + L\kappa_{\xi}) \|\zeta_{T} - \zeta^{*}\| \leq \frac{\varepsilon\kappa_{\zeta}}{1 + \kappa_{\zeta}} \end{aligned}$$

where in the first inequality we use the smoothness Assumption A, and in the second inequality we use (1) in Lemma A.1. As a result,

$$\begin{split} \mathbb{E}[\|\nabla_{\theta}J(\pi_{\theta_{T}})\|] \leq & \frac{\varepsilon}{(\kappa_{\xi}+1)(\kappa_{\zeta}+1)} + \frac{\varepsilon\kappa_{\zeta}}{1+\kappa_{\zeta}} + \varepsilon_{func} + \varepsilon_{reg} + \varepsilon_{data} \\ \leq & \varepsilon + \varepsilon_{func} + \varepsilon_{reg} + \varepsilon_{data} \end{split}$$

C.2 Algorithm Examples

We first introduce a useful assumption:

Assumption G (Diameter). We use Ξ to denote the set of parameters ξ , we assume Ξ is a convex and bounded set with a diameter d > 0.

C.2.1 Example 1: Stochastic Gradient Descent Ascent [20]

Algorithm 2: Direct SGDA 1 Initialize θ_0, ζ_0, ξ_0 2 for t = 0, 1, 2, ...T do 3 Sample $N(s, a, r, s') \sim d^D, a' \sim \pi_{\theta_{t+1}}(s')$ tuples and computing: 4 $\theta_{t+1} \leftarrow \theta_t + \eta_\theta \widehat{\nabla}_\theta \mathcal{L}^D(\theta_t, \zeta_t, \xi_t)$ 5 $\zeta_{t+1} \leftarrow \zeta_t + \eta_\zeta \widehat{\nabla}_\zeta \mathcal{L}^D(\theta_t, \zeta_t, \xi_t)$ 6 $\xi_{t+1} \leftarrow \mathcal{P}_{\xi}(\xi_t - \eta_{\xi} \widehat{\nabla}_{\xi} \mathcal{L}^D(\theta_t, \zeta_t, \xi_t)) // \mathcal{P}_{\xi}$ is the projection operator. 7 end

Adapting from Theorem 4.5 and Proposition 4.11 in [20], we have the following theorem

Theorem C.1. Define $\Delta = \max_{\theta,\zeta} \min_{\xi \in \Xi} \mathcal{L}^D(\theta,\zeta,\xi) - \min_{\xi \in \Xi} \mathcal{L}^D(\theta_0,\zeta_0,\xi)$. Under Assumption A, C, E and G, with step sizes $\eta_{\xi} = \Theta(1/L), \eta_{\zeta} = \eta_{\theta} = \Theta(1/\kappa_{\xi}^2 L)$ and batch size $N = \Theta(\max\{1, \kappa_{\xi}(\kappa_{\xi}+1)^2(\kappa_{\zeta}+1)^2\sigma^2\varepsilon^{-2}\})$, if $T = O(\frac{(\kappa_{\xi}+1)^2(\kappa_{\zeta}+1)^2(\kappa_{\xi}^2L\Delta+\kappa_{\xi}^2L^2D^2)}{\varepsilon^2})$, Algorithm 1 will return us $(\theta_T, \zeta_T, \xi_T)$ satisfying the ε -stationary condition in Eq.(6). In another word, π_{θ_T} satisfies

$$\mathbb{E}[\|J(\pi_{\theta_T})\|] \le \varepsilon + \varepsilon_{data} + \varepsilon_{func} + \varepsilon_{reg}$$
(18)

where ε_{data} is defined in Assumption 2.3, and ε_{func} and ε_{reg} are defined in Theorem 2.7.

C.3 Example 2: Stochastic Recursive Gradient Descent Ascent [8]

In [8], the author presented another algorithm has better dependence on ε . Similarly, we can adapt their algorithm and we ignore the details here.

D Missing details for Algorithm 1

In the following, we will use \mathcal{L}_t^D , \mathcal{L}_t^B and \mathcal{L}_t^{D*} as shortnotes of $\mathcal{L}^D(\theta_t, \zeta_t, \xi_t)$, $\mathcal{L}^B(\theta_t, \zeta_t, \xi_t)$ and $\mathcal{L}^D(\theta_t, \zeta_t^*, \xi_t^*)$, where ζ_t^*, ξ_t^* is the only one saddle point of $\mathcal{L}^D(\theta_t, \zeta, \xi)$. Besides, we use $\nabla_{\theta} \mathcal{L}_t^D$ and $\nabla_{\theta} \mathcal{L}_t^B$ as a shortnote of the gradient averaged over d^D and the gradient averaged over batch, respectively.

Lemma D.1. Suppose we have two empirical gradient estimator $\nabla_{\theta} \mathcal{L}_{t+1}^B$ and $\nabla_{\theta} \mathcal{L}_t^B$ built with the same batch data *B*, under Assumption *A*, we have:

$$\mathbb{E}[\|\nabla_{\theta}\mathcal{L}_{t+1}^{B} - \nabla_{\theta}\mathcal{L}_{t}^{B}\|^{2}]$$

$$\leq 3 \Big(G^{2}L_{w}^{2}C_{Q}^{2}\mathbb{E}[\|\zeta_{t+1} - \zeta_{t}\|^{2}] + G^{2}L_{Q}^{2}C_{W}^{2}\mathbb{E}[\|\xi_{t+1} - \xi_{t}\|^{2}] + H^{2}C_{Q}^{2}C_{W}^{2}\mathbb{E}[\|\theta_{t+1} - \theta_{t}\|^{2}] \Big)$$

Proof.

$$\begin{split} \mathbb{E}[\|\nabla_{\theta}\mathcal{L}_{t+1}^{B} - \nabla_{\theta}\mathcal{L}_{t}^{B}\|^{2}] \\ \leq & \frac{3}{|B|^{2}} \mathbb{E}\Big[B \cdot \sum_{B} \|(1-\gamma)\mathbb{I}[s \in S_{0}] \Big(Q_{t+1}(s,a_{0}) - Q_{t}(s,a_{0})\Big) \nabla_{\theta} \log \pi_{t}(a_{0}|s) \\ & + \gamma w_{t}(s,a) \Big(Q_{t+1}(s',a') - Q_{t}(s',a')\Big) \nabla_{\theta} \log \pi_{t}(a'|s')\|^{2} \\ & + \|(1-\gamma)\mathbb{I}[s \in S_{0}]Q_{t+1}(s,a_{0}) \Big(\nabla_{\theta} \log \pi_{t+1}(a_{0}|s) - \nabla_{\theta} \log \pi_{t}(a_{0}|s)\Big) \\ & + \gamma w_{t}(s,a)Q_{t+1}(s',a') \Big(\nabla_{\theta} \log \pi_{t+1}(a'|s') - \nabla_{\theta} \log \pi_{t}(a'|s')\Big)\|^{2} \\ & + \|\gamma(w_{t+1}(s,a) - w_{t}(s,a))Q_{t+1}(s',a')\nabla_{\theta} \log \pi_{t+1}(a'|s')\|^{2}\Big)\Big] \\ \leq & 3\Big(\gamma^{2}G^{2}L_{w}^{2}C_{Q}^{2}\mathbb{E}[\|\zeta_{t+1} - \zeta_{t}\|^{2}] + G^{2}L_{Q}^{2}\Big((1-\gamma) + \gamma C_{W}\Big)^{2}\mathbb{E}[\|\xi_{t+1} - \xi_{t}\|^{2}] \\ & + H^{2}C_{Q}^{2}\Big((1-\gamma) + \gamma C_{W}\Big)^{2}\mathbb{E}[\|\theta_{t+1} - \theta_{t}\|^{2}]\Big) \\ \leq & 3\Big(G^{2}L_{w}^{2}C_{Q}^{2}\mathbb{E}[\|\zeta_{t+1} - \zeta_{t}\|^{2}] + G^{2}L_{Q}^{2}C_{W}^{2}\mathbb{E}[\|\xi_{t+1} - \xi_{t}\|^{2}] + H^{2}C_{Q}^{2}C_{W}^{2}\mathbb{E}[\|\theta_{t+1} - \theta_{t}\|^{2}]\Big) \end{split}$$

where in the first inequality, we use Young's inequality; in the second one we use Assumption A; in the last one, we use $1 \le C_W$.

Lemma D.2. Under Assumption A, C and D, consider $\pi_{\theta_1}, \pi_{\theta_2}$ parameterized by $\theta_1, \theta_2 \in \Theta$. Θ . Denote (ζ_1^*, ξ_1^*) and (ζ_2^*, ξ_2^*) as the saddle-point of $\max_{\zeta \in Z} \min_{\xi \in \Xi} \mathcal{L}^D(\theta_1, \zeta, \xi)$ and $\max_{\zeta \in Z} \min_{\xi \in \Xi} \mathcal{L}^D(\theta_2, \zeta, \xi)$ respectively, then we have

$$\|\zeta_1^* - \zeta_2^*\| \le \kappa_\mu (\kappa_\xi + 1) \|\theta_1 - \theta_2\| \\ \|\xi_1^* - \xi_2^*\| \le \kappa_\xi (\kappa_\mu + 1) \|\theta_1 - \theta_2\|$$

Proof. With Assumption A and Assumption D, we have

$$|\nabla_{\zeta} \mathcal{L}^{D}(\theta_{2}, \zeta_{1}^{*}, \xi_{1}^{*})|| = ||\nabla_{\zeta} \mathcal{L}^{D}(\theta_{1}, \zeta_{1}^{*}, \xi_{1}^{*}) - \nabla_{\zeta} \mathcal{L}^{D}(\theta_{2}, \zeta_{1}^{*}, \xi_{1}^{*})|| \le L ||\theta_{1} - \theta_{2}||$$
(19)

$$\|\nabla_{\xi}\mathcal{L}^{D}(\theta_{2},\zeta_{1}^{*},\xi_{1}^{*})\| = \|\nabla_{\xi}\mathcal{L}^{D}(\theta_{1},\zeta_{1}^{*},\xi_{1}^{*}) - \nabla_{\xi}\mathcal{L}^{D}(\theta_{2},\zeta_{1}^{*},\xi_{1}^{*})\| \le L\|\theta_{1} - \theta_{2}\|$$
(20)

Recall in Lemma A.1, we know $\Phi_{\theta_2}(\zeta)$ should be a μ_{ζ} -strongly-concave function. Then, we have

$$\begin{split} \|\zeta_{1}^{*} - \zeta_{2}^{*}\| &\leq \frac{1}{\mu_{\zeta}} \|\nabla_{\zeta} \Phi_{\theta_{2}}(\zeta_{1}^{*})\| = \frac{1}{\mu_{\zeta}} \|\nabla_{\zeta} \mathcal{L}^{D}(\theta_{2}, \zeta_{1}^{*}, \phi_{\theta_{2}}(\zeta_{1}^{*}))\| \\ &\leq \frac{1}{\mu_{\zeta}} \|\nabla_{\zeta} \mathcal{L}^{D}(\theta_{2}, \zeta_{1}^{*}, \phi_{\theta_{2}}(\zeta_{1}^{*})) - \nabla_{\zeta} \mathcal{L}^{D}(\theta_{2}, \zeta_{1}^{*}, \xi_{1}^{*})\| + \frac{1}{\mu_{\zeta}} \|\nabla_{\zeta} \mathcal{L}^{D}(\theta_{2}, \zeta_{1}^{*}, \xi_{1}^{*})\| \\ &\leq \frac{1}{\mu_{\zeta}} \|\nabla_{\zeta} \mathcal{L}^{D}(\theta_{2}, \zeta_{1}^{*}, \phi_{\theta_{2}}(\zeta_{1}^{*})) - \nabla_{\zeta} \mathcal{L}^{D}(\theta_{2}, \zeta_{1}^{*}, \xi_{1}^{*})\| + \frac{L}{\mu_{\zeta}} \|\theta_{1} - \theta_{2}\| \\ &\leq \frac{L}{\mu_{\zeta}} \|\phi_{\theta_{2}}(\zeta_{1}^{*}) - \xi_{1}^{*}\| + \frac{L}{\mu_{\zeta}} \|\theta_{1} - \theta_{2}\| \\ &\leq \frac{L}{\mu_{\zeta} \mu_{\xi}} \|\nabla_{\xi} \mathcal{L}^{D}(\theta_{2}, \zeta_{1}^{*}, \xi_{1}^{*})\| + \frac{L}{\mu_{\zeta}} \|\theta_{1} - \theta_{2}\| \\ &\leq \kappa_{\mu} (\kappa_{\xi} + 1) \|\theta_{1} - \theta_{2}\| \end{split}$$

where in the first step, we use Lemma A.3; in the fourth inequality, we use Assumption A; in the fifth inequality, we use the Assumption C that, given $\theta_2, \zeta_1^*, \mathcal{L}^D(\theta_2, \zeta_1^*, \xi)$ is μ_{ξ} -strongly-convex w.r.t. ξ and $\phi_{\theta_2}(\zeta_1^*)$ is the optimum of it; in the last inequality, we use Eq.(19) again.

We can give a similarly discussion for $\|\xi_1^* - \xi_2^*\|$:

$$\begin{split} \|\xi_{1}^{*} - \xi_{2}^{*}\| &\leq \frac{1}{\mu_{\xi}} \|\nabla_{\xi} \Psi_{\theta_{2}}(\xi_{1}^{*})\| = \frac{1}{\mu_{\xi}} \|\nabla_{\xi} \mathcal{L}^{D}(\theta_{2}, \psi_{\theta_{2}}(\xi_{1}^{*}), \xi_{1}^{*})\| \\ &\leq \frac{1}{\mu_{\xi}} \|\nabla_{\xi} \mathcal{L}^{D}(\theta_{2}, \psi_{\theta_{2}}(\xi_{1}^{*}), \xi_{1}^{*}) - \nabla_{\xi} \mathcal{L}^{D}(\theta_{2}, \zeta_{1}^{*}, \xi_{1}^{*})\| + \frac{1}{\mu_{\xi}} \|\nabla_{\xi} \mathcal{L}^{D}(\theta_{2}, \zeta_{1}^{*}, \xi_{1}^{*}))\| \\ &\leq \frac{1}{\mu_{\xi}} \|\nabla_{\xi} \mathcal{L}^{D}(\theta_{2}, \psi_{\theta_{2}}(\xi_{1}^{*}), \xi_{1}^{*}) - \nabla_{\xi} \mathcal{L}^{D}(\theta_{2}, \zeta_{1}^{*}, \xi_{1}^{*})\| + \frac{L}{\mu_{\xi}} \|\theta_{1} - \theta_{2}\| \\ &\leq \frac{L}{\mu_{\xi}} \|\zeta_{1}^{*} - \psi_{\theta_{2}}(\xi_{1}^{*})\| + \frac{L}{\mu_{\xi}} \|\theta_{1} - \theta_{2}\| \\ &\leq \frac{L}{\mu_{\xi}\mu_{\zeta}} \|\nabla_{\zeta} \mathcal{L}^{D}(\theta_{2}, \zeta_{1}^{*}, \xi_{1}^{*})\| + \frac{L}{\mu_{\xi}} \|\theta_{1} - \theta_{2}\| \\ &\leq \kappa_{\xi}(\kappa_{\mu} + 1)\|\theta_{1} - \theta_{2}\| \end{split}$$

Lemma D.3 (Relate the shift of ζ_t and ξ_t with θ_t). We consider the Assumptions A, C, F and D. Denote $(\theta_t, \zeta_t, \xi_t)$ as the parameter value at the beginning at the step t in Algorithm 1, and denote $(\zeta_t^*, \xi_t^*) \in Z \times \Xi$ as the only saddle point for $\mathcal{L}^D(\theta_t, \zeta, \xi)$ given θ_t . Recall the Oracle in Definition 4.1 that, for arbitrary t iteration, it will return us ζ_{t+1}, ξ_{t+1} satisfying

$$\mathbb{E}[\|\zeta_{t+1} - \zeta_{t+1}^*\|^2 + \|\xi_{t+1} - \xi_{t+1}^*\|^2] \le \frac{\beta}{2}\mathbb{E}[\|\zeta_t - \zeta_{t+1}^*\|^2 + \|\xi_t - \xi_{t+1}^*\|^2] + c$$

where $0 < \beta/2 \leq 1$. Then, we have:

$$\mathbb{E}[\|\zeta_{t+1} - \zeta_t\|^2 + \|\xi_{t+1} - \xi_t\|^2] \le 6\beta^{t+1}d^2 + 6\eta_\theta^2 C_{\zeta,\mu} \sum_{\tau=0}^t \beta^{t-\tau} \mathbb{E}[\|g_\theta^\tau\|^2] + \frac{6c}{1-\beta}$$

where d is the diameter defined in Assumption F, and $C_{\zeta,\mu}$ is a short note of $\kappa^2_{\mu}(\kappa_{\xi}+1)^2 + \kappa^2_{\xi}(\kappa_{\mu}+1)^2$.

Proof. We will use $\Delta_t(\zeta, \xi)$ to denote $\mathbb{E}[\|\zeta - \zeta_t^*\|^2 + \|\xi - \xi_t^*\|^2]$. We first study some useful properties of $\Delta_t(\zeta, \xi)$.

Property 1 For $t \ge 1$

$$\begin{split} \Delta_t(\zeta_{t-1}^*, \xi_{t-1}^*) = & \mathbb{E}[\|\zeta_t^* - \zeta_{t-1}^*\|^2 + \|\xi_t^* - \xi_{t-1}^*\|^2] \\ \leq & C_{\zeta,\mu} \mathbb{E}[\|\theta_t - \theta_{t-1}\|^2] \\ = & \eta_\theta^2 C_{\zeta,\mu} \mathbb{E}[\|g_\theta^{t-1}\|^2] \end{split}$$

where in the inequality, we use Lemma D.2; and the last equality results from the update rule $\theta_t = \theta_{t-1} + \eta_{\theta} g_{\theta}^{t-1}$

$$\begin{aligned} \mathbf{Property 2} \quad & \text{For } t \geq 0, \\ \Delta_t(\zeta_t, \xi_t) \leq & \frac{\beta}{2} \Delta_t(\zeta_{t-1}, \xi_{t-1}) + c = \frac{\beta}{2} \mathbb{E}[\|\zeta_{t-1} - \zeta_t^*\|^2 + \|\xi_{t-1} - \xi_t^*\|^2] + c \\ & \leq & \beta \mathbb{E}[\|\zeta_{t-1} - \zeta_{t-1}^*\|^2 + \|\xi_{t-1} - \xi_{t-1}^*\|^2 + \|\zeta_t^* - \zeta_{t-1}^*\|^2 + \|\xi_t^* - \xi_{t-1}^*\|^2] + c \\ & = & \beta \Delta_{t-1}(\zeta_{t-1}, \xi_{t-1}) + \beta \Delta_t(\zeta_{t-1}^*, \xi_{t-1}^*) + c \\ & \leq & \beta^t \Delta_0(\zeta_0, \xi_0) + \sum_{\tau=1}^t \beta^{t-\tau+1} \Delta_\tau(\zeta_{\tau-1}^*, \xi_{\tau-1}^*) + \sum_{\tau=0}^{t-1} \beta^\tau c \\ & \leq & \beta^{t+1} d^2 + \eta_\theta^2 C_{\zeta, \mu} \sum_{\tau=0}^{t-1} \beta^{t-\tau} \mathbb{E}[\|g_\theta^\tau\|^2] + \sum_{\tau=0}^t \beta^\tau c \\ & \leq & \beta^{t+1} d^2 + \eta_\theta^2 C_{\zeta, \mu} \sum_{\tau=0}^{t-1} \beta^{t-\tau} \mathbb{E}[\|g_\theta^\tau\|^2] + \frac{c}{1-\beta} \end{aligned}$$

where the first inequality is because of the property of the Oracle; for the second inequality we use Young's inequality; In the last step, we use

$$\Delta_0(\zeta_0,\xi_0) = \mathbb{E}[\|\zeta_0 - \zeta_0^*\|^2 + \|\xi_0 - \xi_0^*\|^2] \le \frac{\beta}{2} \mathbb{E}[\|\zeta_{-1} - \zeta_0^*\|^2 + \|\xi_{-1} - \xi_0^*\|^2] + c \le \beta d^2 + c$$

With the two properties above, we can bound:

$$\begin{split} & \mathbb{E}[\|\zeta_{t+1} - \zeta_t\|^2 + \|\xi_{t+1} - \xi_t\|^2] \\ \leq 3\mathbb{E}[\|\zeta_{t+1} - \zeta_{t+1}^*\|^2 + \|\xi_{t+1} - \xi_{t+1}^*\|^2 + \|\zeta_{t+1}^* - \zeta_t^*\|^2 + \|\xi_{t+1}^* - \xi_t^*\|^2 + \|\zeta_t^* - \zeta_t\|^2 + \|\xi_t^* - \xi_t\|^2] \\ = 3\Delta_{t+1}(\zeta_{t+1}, \xi_{t+1}) + 3\Delta_{t+1}(\zeta_t^*, \xi_t^*) + 3\Delta_t(\zeta_t, \xi_t) \\ \leq 3\beta^{t+2}d^2 + 3\eta_\theta^2 C_{\zeta,\mu} \sum_{\tau=0}^t \beta^{t-\tau+1} \mathbb{E}[\|g_\theta^\tau\|^2] + 3\eta_\theta^2 C_{\zeta,\mu} \mathbb{E}[\|g_\theta^t\|^2] + 3\beta^{t+1}d^2 + 3\eta_\theta^2 C_{\zeta,\mu} \sum_{\tau=0}^{t-1} \beta^{t-\tau} \mathbb{E}[\|g_\theta^\tau\|^2] \\ &+ \frac{6c}{1-\beta} \\ = 3(1+\beta)\beta^{t+1}d^2 + 3\eta_\theta^2 C_{\zeta,\mu} \sum_{\tau=0}^t (1+\beta)\beta^{t-\tau} \mathbb{E}[\|g_\theta^\tau\|^2] + \frac{6c}{1-\beta} \\ \leq 6\beta^{t+1}d^2 + 6\eta_\theta^2 C_{\zeta,\mu} \sum_{\tau=0}^t \beta^{t-\tau} \mathbb{E}[\|g_\theta^\tau\|^2] + \frac{6c}{1-\beta} \end{split}$$

where for the first one we use an extended version of Young's inequality $\|\sum_{i=1}^{k} x_i\|^2 \le k \sum_{i=1}^{k} \|x_i\|^2$; in the second inequality, we use the Property 1 and 2 to give the upper bound; in the third inequality, we use the fact that $0 < \beta \le 1$.

Lemma D.4. Under the same condition of Lemma D.3 above, with an additional constraint $\beta \le (1-\alpha)^2/2$ and an additional Assumption E, for $t \ge 0$, we have:

$$\begin{split} & \mathbb{E}[\|g_{\theta}^{t+1} - \nabla_{\theta} J(\theta_{t+1})\|^{2}] \\ \leq & 3(\varepsilon_{data} + \varepsilon_{func} + \varepsilon_{reg})^{2} + 3(1-\alpha)^{2t+2} \mathbb{E}[\|g_{\theta}^{0} - \nabla_{\theta} \mathcal{L}_{0}^{D}\|^{2}] + \frac{6\alpha\sigma^{2}}{|B|} \\ & + \left(6L^{2}(\beta^{t+2}d^{2} + \frac{c}{1-\beta}) + 108C_{w,Q}\left(\frac{\beta(1-\alpha)^{2(t+2)}}{(1-\alpha)^{2}-\beta}d^{2} + \frac{c}{\alpha(2-\alpha)(1-\beta)}\right) \right) \\ & + \sum_{i=0}^{t} \left(108\eta_{\theta}^{2}(1-\alpha)^{2(t-i+1)}\left(2C_{\zeta,\mu}C_{w,Q} + H^{2}C_{Q}^{2}C_{W}^{2}\right) + 6L^{2}\eta_{\theta}^{2}C_{\zeta,\mu}\beta^{t-i+1}\right) \mathbb{E}[\|g_{\theta}^{i}\|^{2}] \end{split}$$

where $\varepsilon_{data}, \varepsilon_{func}, \varepsilon_{reg}$ are the same as those in Theorem 2.7, and

$$C_{w,Q} := G^2 L_w^2 C_Q^2 + G^2 L_Q^2 C_W^2$$

Proof. Recall that we will use $\nabla_{\theta} \mathcal{L}_{t}^{B}$, $\nabla_{\theta} \mathcal{L}_{t}^{D}$ and $\nabla_{\theta} \mathcal{L}_{t}^{D*}$ as a shortnote of $\nabla_{\theta} \mathcal{L}^{B}(\theta_{t}, \zeta_{t}, \xi_{t})$, $\nabla_{\theta} \mathcal{L}^{D}(\theta_{t}, \zeta_{t}^{*}, \xi_{t}^{*})$ respectively. First we can use the Young's inequality to obtain

$$\mathbb{E}[\|g_{\theta}^{t+1} - \nabla_{\theta}J(\theta_{t+1})\|^{2}]$$

$$\leq 3 \underbrace{\mathbb{E}[\|\nabla_{\theta}\mathcal{L}_{t+1}^{D*} - \nabla_{\theta}J(\theta_{t+1})\|^{2}]}_{Bias \ (Bounded \ in \ Theorem \ 2.7)} + 3 \underbrace{\mathbb{E}[\|g_{\theta}^{t+1} - \nabla_{\theta}\mathcal{L}_{t+1}^{D}\|^{2}]}_{p_{1}} + 3 \underbrace{\mathbb{E}[\|\nabla_{\theta}\mathcal{L}_{t+1}^{D} - \nabla_{\theta}\mathcal{L}_{t+1}^{D*}\|^{2}]}_{p_{2}}$$

Since the first term has already been bounded in Theorem 2.7. Next, we bound p_1 and p_2 :

Upper bound p_1 We again use $C_{\zeta,\xi}$ as a short note of $\kappa^2_{\mu}(\kappa_{\xi}+1)^2 + \kappa^2_{\xi}(\kappa_{\mu}+1)^2$. From Lemma D.3, we know that,

$$\mathbb{E}[\|\zeta_{t+1} - \zeta_t\|^2 + \|\xi_{t+1} - \xi_t\|^2] \le 6\beta^{t+1}d^2 + 6\eta_\theta^2 C_{\zeta,\mu} \sum_{\tau=0}^t \beta^{t-\tau} \mathbb{E}[\|g_\theta^\tau\|^2] + \frac{6c}{1-\beta}$$

Then, we have

$$\begin{aligned} &\leq (1-\alpha)^{2} \mathbb{E}[\|g_{\theta}^{t} - \nabla_{\theta} \mathcal{L}_{t}^{D}\|^{2}] + \frac{2\alpha^{2}\sigma^{2}}{|B|} \\ &+ 6(1-\alpha)^{2} \left(G^{2} L_{w}^{2} C_{Q}^{2} \mathbb{E}[\|\zeta_{t+1} - \zeta_{t}\|^{2}] + G^{2} L_{Q}^{2} C_{W}^{2} \mathbb{E}[\|\xi_{t+1} - \xi_{t}\|^{2}] + H^{2} C_{Q}^{2} C_{W}^{2} \mathbb{E}[\|\theta_{t+1} - \theta_{t}\|^{2}] \right) \\ &\leq (1-\alpha)^{2t+2} \mathbb{E}[\|g_{\theta}^{0} - \nabla_{\theta} \mathcal{L}_{0}^{D}\|^{2}] + \frac{2\alpha^{2}\sigma^{2}}{|B|} \frac{1-(1-\alpha)^{2t+2}}{1-(1-\alpha)^{2}} \\ &+ 6\mathbb{E}\Big[\sum_{i=0}^{t} (1-\alpha)^{2(t-i+1)} \left(G^{2} L_{w}^{2} C_{Q}^{2} \|\zeta_{i+1} - \zeta_{i}\|^{2} + G^{2} L_{Q}^{2} C_{W}^{2} \|\xi_{i+1} - \xi_{i}\|^{2} + H^{2} C_{Q}^{2} C_{W}^{2} \|\theta_{i+1} - \theta_{i}\|^{2} \right) \\ &\leq (1-\alpha)^{2t+2} \mathbb{E}[\|g_{\theta}^{0} - \nabla_{\theta} \mathcal{L}_{0}^{D}\|^{2}] + \frac{2\alpha\sigma^{2}}{|B|} + 36 \sum_{i=0}^{t} (1-\alpha)^{2(t-i+1)} C_{w,Q} (\beta^{i+1}d^{2} + \frac{c}{1-\beta}) \\ &\quad (\alpha < 1) \\ &+ 36\eta_{\theta}^{2} \sum_{i=0}^{t} \left(C_{\zeta,\mu} C_{w,Q} \sum_{\tau=i}^{t} (1-\alpha)^{2(t-\tau+1)} \beta^{\tau-i} + (1-\alpha)^{2(t-i+1)} H^{2} C_{Q}^{2} C_{W}^{2} \right) \mathbb{E}[\|g_{\theta}^{i}\|^{2}] \\ &\leq (1-\alpha)^{2t+2} \mathbb{E}[\|g_{\theta}^{0} - \nabla_{\theta} \mathcal{L}_{0}^{D}\|^{2}] + \frac{2\alpha\sigma^{2}}{|B|} + 36C_{w,Q} \left(\frac{\beta(1-\alpha)^{2(t-2)}}{(1-\alpha)^{2}-\beta} d^{2} + \frac{(1-\alpha)^{2c}}{(1-(1-\alpha)^{2})(1-\beta)} \right) \\ &+ 36\eta_{\theta}^{2} \sum_{i=0}^{t} (1-\alpha)^{2(t-i+1)} \left(C_{\zeta,\mu} C_{w,Q} \frac{(1-\alpha)^{2}}{|B|} + 36C_{w,Q} \left(\frac{\beta(1-\alpha)^{2(t+2)}}{(1-\alpha)^{2}-\beta} d^{2} + \frac{(1-\alpha)^{2c}}{(1-(1-\alpha)^{2})(1-\beta)} \right) \\ &+ 36\eta_{\theta}^{2} \sum_{i=0}^{t} (1-\alpha)^{2(t-i+1)} \left(C_{\zeta,\mu} C_{w,Q} + H^{2} C_{Q}^{2} C_{W}^{2} \right) \mathbb{E}[\|g_{\theta}^{i}\|^{2}] \end{aligned}$$

where the fourth equality because $\mathbb{E}[\nabla_{\theta}\mathcal{L}_{t}^{B}] = \nabla_{\theta}\mathcal{L}_{t}^{D}$ holds for all t and so the cross terms has 0 expectation; the first inequality is because variance is less than the second momentum; the second inequality we apply Lemma D.1 and Assumption A; in the last but two inequality, we apply the summation formula of equal ratio sequence and use the fact that $0 < \alpha \le 1, \beta \le 1$; in the last step, we use our condition $\beta \le (1 - \alpha)^2/2$

Upper bound p_2 Next, we give an upper bound for p_2 . From the Property 2 in Lemma D.3, we know that

$$\begin{aligned} \Delta_{t+1}(\zeta_{t+1}, \xi_{t+1}) = & \mathbb{E}[\|\zeta_{t+1} - \zeta_{t+1}^*\|^2] + \mathbb{E}[\|\xi_{t+1} - \xi_{t+1}^*\|^2] \\ \leq & \beta^{t+2} d^2 + \eta_{\theta}^2 C_{\zeta,\mu} \sum_{\tau=0}^t \beta^{t-\tau+1} \mathbb{E}[\|g_{\theta}^{\tau}\|^2] + \frac{c}{1-\beta} \end{aligned}$$

As a result

$$p_{2} = \mathbb{E}[\|\nabla_{\theta}\mathcal{L}_{t+1}^{D} - \nabla_{\theta}\mathcal{L}_{t+1}^{D*}\|^{2}] \leq 2L^{2}\mathbb{E}[\|\zeta_{t+1} - \zeta_{t+1}^{*}\|^{2} + \|\xi_{t+1} - \xi_{t+1}^{*}\|^{2}]$$
$$\leq 2L^{2}\Big(\beta^{t+2}d^{2} + \eta_{\theta}^{2}C_{\zeta,\mu}\sum_{\tau=0}^{t}\beta^{t-\tau+1}\mathbb{E}[\|g_{\theta}^{\tau}\|^{2}] + \frac{c}{1-\beta})$$

Combine these two results we can finish the proof:

$$\begin{split} \mathbb{E}[\|g_{\theta}^{t+1} - \nabla_{\theta} J(\theta_{t+1})\|^{2}] &\leq 3\mathbb{E}[\|\nabla_{\theta} \mathcal{L}_{t+1}^{D*} - \nabla_{\theta} J(\theta_{t+1})\|^{2}] + 3p_{1} + 3p_{2} \\ &\leq 3(\varepsilon_{data} + \varepsilon_{func} + \varepsilon_{reg})^{2} + 3(1-\alpha)^{2t+2}\mathbb{E}[\|g_{\theta}^{0} - \nabla_{\theta} \mathcal{L}_{0}^{D}\|^{2}] + \frac{6\alpha\sigma^{2}}{|B|} \\ &+ 108C_{w,Q} = \left(\frac{\beta(1-\alpha)^{2(t+2)}}{(1-\alpha)^{2} - \beta} d^{2} + \frac{(1-\alpha)^{2}c}{(1-(1-\alpha)^{2})(1-\beta)}\right) \\ &+ 108\eta_{\theta}^{2} \sum_{i=0}^{t} (1-\alpha)^{2(t-i+1)} \left(2C_{\zeta,\mu}C_{w,Q} + H^{2}C_{Q}^{2}C_{W}^{2}\right)\mathbb{E}[\|g_{\theta}^{i}\|^{2}] \\ &+ 6L^{2} \left(\beta^{t+2}d^{2} + \eta_{\theta}^{2}C_{\zeta,\mu}\sum_{\tau=0}^{t} \beta^{t-\tau+1}\mathbb{E}[\|g_{\theta}^{\tau}\|^{2}] + \frac{c}{1-\beta}\right) \\ &\leq 3(\varepsilon_{data} + \varepsilon_{func} + \varepsilon_{reg})^{2} + 3(1-\alpha)^{2t+2}\mathbb{E}[\|g_{\theta}^{0} - \nabla_{\theta}\mathcal{L}_{0}^{D}\|^{2}] + \frac{6\alpha\sigma^{2}}{|B|} \\ &+ \left(6L^{2}(\beta^{t+2}d^{2} + \frac{c}{1-\beta}) + 108C_{w,Q} \left(\frac{\beta(1-\alpha)^{2(t+2)}}{(1-\alpha)^{2} - \beta}d^{2} + \frac{c}{\alpha(2-\alpha)(1-\beta)}\right) \right) \\ &+ \sum_{i=0}^{t} \left(108\eta_{\theta}^{2}(1-\alpha)^{2(t-i+1)} \left(2C_{\zeta,\mu}C_{w,Q} + H^{2}C_{Q}^{2}C_{W}^{2}\right) + 6L^{2}\eta_{\theta}^{2}C_{\zeta,\mu}\beta^{t-i+1}\right)\mathbb{E}[\|g_{\theta}^{i}\|^{2}] \\ & \square \end{split}$$

Proposition 4.2. Under Assumption A, $J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}, s_0 \sim \nu_0} [\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)]$ is L_J smooth with

$$L_J := \frac{H}{(1-\gamma)^2} + \frac{(1+\gamma)G^2}{(1-\gamma)^3}$$

Proof. Recall that,

$$\nabla_{\theta} J(\pi) = \int_{\tau} p(\tau|\theta) \sum_{i=0}^{\infty} \gamma^{i} r_{i} \sum_{j=0}^{i} \nabla_{\theta} \log \pi_{\theta}(a_{j}|s_{j}) d\tau$$

Therefore,

$$\begin{aligned} \nabla_{\theta}^{2} J(\pi) &= \int_{\tau} p(\tau|\theta) \sum_{i=0}^{\infty} \gamma^{i} r_{i} \sum_{j=0}^{i} \nabla_{\theta}^{2} \log \pi_{\theta}(a_{j}|s_{j}) d\tau \\ &+ \int_{\tau} p(\tau|\theta) \nabla_{\theta} \log p(\tau|\theta) \sum_{i=0}^{\infty} \gamma^{i} r_{i} \sum_{j=0}^{i} \nabla_{\theta} \log \pi_{\theta}(a_{j}|s_{j}) d\tau \\ &= \int_{\tau} p(\tau|\theta) \sum_{i=0}^{\infty} \gamma^{i} r_{i} \sum_{j=0}^{i} \nabla_{\theta}^{2} \log \pi_{\theta}(a_{j}|s_{j}) d\tau \\ &+ \int_{\tau} p(\tau|\theta) \sum_{i=0}^{\infty} \gamma^{i} r_{i} \left(\sum_{j=0}^{i} \nabla_{\theta} \log \pi(a_{t}|s_{t})\right) \left(\sum_{j=0}^{i} \nabla_{\theta} \log \pi(a_{t}|s_{t})\right)^{\top} d\tau \end{aligned}$$

Therefore,

$$\begin{split} \|\nabla_{\theta}^{2}J(\pi)\|_{op} &\leq \int_{\tau} p(\tau|\theta) \sum_{i=0}^{\infty} \gamma^{i} \sum_{j=0}^{i} \|\nabla_{\theta}^{2} \log \pi_{\theta}(a_{j}|s_{j})\|_{op} d\tau \\ &+ \int_{\tau} p(\tau|\theta) \sum_{i=0}^{\infty} \gamma^{i} \| \left(\sum_{j=0}^{i} \nabla_{\theta} \log \pi(a_{t}|s_{t})\right) \left(\sum_{j=0}^{i} \nabla_{\theta} \log \pi(a_{t}|s_{t})\right)^{\top} \|_{op} d\tau \\ &\leq \sum_{i=0}^{\infty} \gamma^{i} (i+1)H + \sum_{i=0}^{\infty} \gamma^{i} (i+1)^{2} G^{2} \\ &= \frac{H}{(1-\gamma)^{2}} + \frac{(1+\gamma)G^{2}}{(1-\gamma)^{3}} \end{split}$$

Theorem 4.3. Under Assumption A-F and H, given arbitrary ε , by choosing Algorithm 3 as the Oracle, Algorithm 1 will return us a policy $\pi_{\theta_{out}}$, satisfying

$$\mathbb{E}[\|\nabla_{\theta} J(\pi_{\theta_T})\|] \le \varepsilon + \sqrt{3}(\varepsilon_{reg} + \varepsilon_{data} + \varepsilon_{func})$$

if the hyper-parameters in Alg. 1 and 3 satisfy the following constraints:

$$\begin{split} T = & [\max\{96, \frac{16L_J}{\varepsilon^2}, \frac{16}{(1-\gamma)\varepsilon^2} \sqrt{120\left(2C_{\zeta,\mu}C_{w,Q} + H^2C_Q^2C_W^2\right)}, \frac{864C_{w,Q}d^2}{\varepsilon^2}\}] = O(\varepsilon^{-2});\\ |B| = & [\max\{1, \frac{12\sigma^2}{\varepsilon^2}\}]; \quad |N| = [\frac{96(L^2 + 20C_{w,Q})\sigma^2}{\min\{\frac{\mu_{\zeta}\eta_{\zeta}}{4}, \frac{\mu_{\xi}\eta_{\xi}}{4}\}\varepsilon^2}(\frac{\eta_{\zeta}}{\mu_{\zeta}} + \frac{\eta_{\xi}}{\mu_{\xi}})]; \quad K = c_{oracle}\log(\frac{1}{\beta});\\ \alpha = & 0.9; \quad \beta = \min\{\frac{\varepsilon^2}{L^2}, \frac{(1-\gamma)^2\varepsilon^4}{C_{\zeta,\mu}L^2}, \frac{\alpha}{2}(1-\alpha)^2\}; \quad B_0 = [\frac{8\sigma^2}{\varepsilon^2}]\\ \eta_{\theta} = \min\{\frac{1}{2L_J}, \left(\left[\frac{C_{\zeta,\mu}L^2\beta}{6(1-\beta)} + 120\left(2C_{\zeta,\mu}C_{w,Q} + H^2C_W^2C_Q^2\right)\right]\right)^{-1/2}\} \end{split}$$

where $[\cdot]$ is the upper rounding function, $C_{w,Q} = G^2 L_w^2 C_Q^2 + G^2 C_W^2 L_Q^2$, $C_{\zeta,\mu} = \kappa_{\mu}^2 (\kappa_{\xi} + 1)^2 + \kappa_{\xi}^2 (\kappa_{\mu} + 1)^2$, L_J is defined in Prop. 4.2, η_{ζ} and η_{ξ} satisfy the constraints in Theorem E.1 and c_{oracle} is an independent constant.

Besides, the total gradient computation to obtain θ_{out} should be $|B_0| + |B| \cdot T + |N| \cdot K \cdot T = O(\varepsilon^{-4})$.

Proof.

$$J(\theta_{T+1}) = J(\theta_T + \eta_\theta g_\theta^T)$$

$$\geq J(\theta_T) + \eta_\theta (g_\theta^T)^\top \nabla_\theta J(\theta_T) - \frac{\eta_\theta^2 L_J}{2} \|g_\theta^T\|^2$$

$$= J(\theta_T) + \frac{\eta_\theta}{2} \|\nabla_\theta J(\theta_T)\|^2 - \frac{\eta_\theta}{2} \|g_\theta^T - \nabla_\theta J(\theta_T)\|^2 + (\frac{\eta_\theta}{2} - \frac{\eta_\theta^2 L_J}{2}) \|g_\theta^T\|^2$$

$$\geq J(\theta_T) + \frac{\eta_\theta}{2} \|\nabla_\theta J(\theta_T)\|^2 - \frac{\eta_\theta}{2} \|g_\theta^T - \nabla_\theta J(\theta_T)\|^2 + \frac{\eta_\theta}{4} \|g_\theta^T\|^2$$

$$\geq J(\theta_0) + \frac{\eta_\theta}{2} \sum_{t=0}^T \|\nabla_\theta J(\theta_t)\|^2 - \frac{\eta_\theta}{2} \Big(\sum_{t=0}^T \|g_\theta^t - \nabla_\theta J(\theta_t)\|^2 - \frac{1}{2} \|g_\theta^t\|^2 \Big)$$

where in the second equation, we use the fact that $(g_{\theta}^T)^{\top} \nabla_{\theta} J(\theta_T) = \frac{1}{2} \|\nabla_{\theta} J(\theta_T)\|^2 + \frac{1}{2} \|g_{\theta}^T\|^2 - \frac{1}{2} \|g_{\theta}^T - \nabla_{\theta} J(\theta_T)\|^2$; in the second inequality, we add a constraint for η_{θ} that $\eta_{\theta} \leq \frac{1}{2L_J}$; Next, we give a upper bound for p with Lemma D.4:

$$\begin{split} p &= \sum_{t=0}^{T} \|g_{\theta}^{\tau} - \nabla_{\theta} J(\theta_{t})\|^{2} - \frac{1}{2} \|g_{\theta}^{t}\|^{2} \\ &\leq \sum_{t=0}^{T} \left\{ 3(\varepsilon_{data} + \varepsilon_{func} + \varepsilon_{reg})^{2} + 3(1-\alpha)^{2t+2} \mathbb{E}[\|g_{\theta}^{0} - \nabla_{\theta} \mathcal{L}_{0}^{D}\|^{2}] + \frac{6\alpha\sigma^{2}}{|B|} \\ &+ \left(6L^{2}(\beta^{t+2}d^{2} + \frac{c}{1-\beta}) + 108C_{w,Q} \left(\frac{\beta(1-\alpha)^{2(t+2)}}{(1-\alpha)^{2}-\beta} d^{2} + \frac{c}{\alpha(2-\alpha)(1-\beta)} \right) \right. \\ &+ \sum_{i=0}^{t} \left(108\eta_{\theta}^{2}(1-\alpha)^{2(t-i+1)} \left(2C_{\zeta,\mu}C_{w,Q} + H^{2}C_{Q}^{2}C_{w}^{2} \right) + 6L^{2}\eta_{\theta}^{2}C_{\zeta,\mu}\beta^{t-i+1} \right) \mathbb{E}[\|g_{\theta}^{i}\|^{2}] - \frac{1}{2}\mathbb{E}[\|g_{\theta}^{i}\|^{2}] \right\} \\ &\leq 3T(\varepsilon_{data} + \varepsilon_{func} + \varepsilon_{reg})^{2} + \left(\frac{6L^{2}}{1-\beta} + \frac{108C_{w,Q}}{\alpha(2-\alpha)(1-\beta)} \right) Tc \\ &+ \frac{3}{1-(1-\alpha)^{2}} \mathbb{E}[\|g_{\theta}^{0} - \nabla_{\theta}\mathcal{L}_{0}^{D}\|^{2}] + \frac{6\alphaT\sigma^{2}}{|B|} + \left(\frac{6\betaL^{2}}{1-\beta} + \frac{108\beta(1-\alpha)^{2}C_{w,Q}}{(1-(1-\alpha)^{2})((1-\alpha)^{2}-\beta)} \right) \right) d^{2} \\ &+ \sum_{t=0}^{T} \mathbb{E}[\|g_{\theta}^{i}\|^{2}] \left\{ -\frac{1}{2} + 108\eta_{\theta}^{2} \sum_{i=1}^{T-t+1} \left[\frac{C_{\zeta,\mu}L^{2}\beta^{i}}{18} + (1-\alpha)^{2i} \left(2C_{\zeta,\mu}C_{w,Q} + H^{2}C_{Q}^{2}C_{w}^{2} \right) \right] \right\} \\ &\leq 3T(\varepsilon_{data} + \varepsilon_{func} + \varepsilon_{reg})^{2} + \left(\frac{6L^{2}}{1-\beta} + \frac{108C_{w,Q}}{\alpha(2-\alpha)(1-\beta)} \right) Tc \\ &+ \frac{3}{\alpha} \mathbb{E}[\|g_{\theta}^{0} - \nabla_{\theta}\mathcal{L}_{0}^{D}\|^{2}] + \frac{6\alphaT\sigma^{2}}{|B|} + \left(\frac{6\betaL^{2}}{1-\beta} + 108C_{w,Q} \frac{\beta}{\alpha((1-\alpha)^{2}-\beta)} \right) \right) d^{2} \\ &+ \sum_{t=0}^{T} \mathbb{E}[\|g_{\theta}^{i}\|^{2}] \left(-\frac{1}{2} + 108\eta_{\theta}^{2} \left[\frac{C_{\zeta,\mu}L^{2}\beta}{18(1-\beta)} + \frac{(1-\alpha)^{2}}{1-(1-\alpha)^{2}} \left(2C_{\zeta,\mu}C_{w,Q} + H^{2}C_{Q}^{2}C_{w}^{2} \right) \right] \right\} \\ &\leq 3T(\varepsilon_{data} + \varepsilon_{func} + \varepsilon_{reg})^{2} + \left(\frac{6L^{2}}{1-\beta} + \frac{108C_{w,Q}}{\alpha(2-\alpha)(1-\beta)} \right) Tc \\ &+ \frac{3}{\alpha} \mathbb{E}[\|g_{\theta}^{0} - \nabla_{\theta}\mathcal{L}_{0}^{0}\|^{2}] + \frac{6\alphaT\sigma^{2}}{|B|} + \left(\frac{6\betaL^{2}}{1-\beta} + 108C_{w,Q} \frac{\alpha((1-\alpha)^{2}-2}{\alpha((1-\alpha)^{2}-(1-\alpha)^{2}/2)} \right) \right) d^{2} \\ &+ \sum_{t=0}^{T} \mathbb{E}[\|g_{\theta}^{i}\|^{2}] \left(-\frac{1}{2} + 108\eta_{\theta}^{2} \left[\frac{C_{\zeta,\mu}L^{2}\beta}{1-\beta} + \frac{108C_{w,Q}}{\alpha(2-\alpha)(1-\beta)} \right) Tc \\ &+ \frac{3}{\alpha} \mathbb{E}[\|g_{\theta}^{i}\|^{2}\right] \left(-\frac{1}{2} + 108\eta_{\theta}^{2} \left[\frac{C_{\zeta,\mu}L^{2}\beta}{1-\beta} + \frac{108C_{w,Q}}{\alpha((1-\alpha)^{2}-(1-\alpha)^{2}/2)} \right) \right) d^{2} \\ &+ \sum_{t=0}^{T} \mathbb{E}[\|g_{\theta}^{i}\|^{2}\right] \left(-\frac{1}{2} + 108\eta_{\theta}^{2} \left[\frac{C_{\zeta,\mu}L^{2}\beta}{1-\beta} + \frac{108C_{w,Q}}{\alpha((1-\alpha)^{2}-(1-\alpha)^{2}/$$

$$+\frac{3}{\alpha}\mathbb{E}[\|g_{\theta}^{0}-\nabla_{\theta}\mathcal{L}_{0}^{D}\|^{2}]+\frac{6\alpha T\sigma^{2}}{|B|}+(\frac{6\beta L^{2}}{1-\beta}+108C_{w,Q})d^{2}$$

In the first, second and third inequality, we use the fact that $0 < (1-\alpha) < 1, 0 < \beta \le \alpha (1-\alpha)^2/2 \le (1-\alpha)^2/2$. In the fourth inequality, we add the following constraint to drop the terms containing $||g_{\theta}||$:

$$\eta_{\theta} \le \left(108 \left[\frac{C_{\zeta,\mu} L^2 \beta}{18(1-\beta)} + \frac{1}{\alpha} \left(2C_{\zeta,\mu} C_{w,Q} + H^2 C_{\mathcal{Q}}^2 C_{\mathcal{W}}^2 \right) \right] \right)^{-1/2}$$
(22)

Therefore,

$$\begin{split} \frac{1}{T+1} \sum_{t=0}^{T} \|\nabla_{\theta} J(\theta_{\tau})\|^{2} &\leq \frac{2}{(T+1)\eta_{\theta}} (J(\theta_{T}) - J(\theta_{0})) + \frac{1}{T+1} \sum_{\tau=0}^{T} \left(\|g_{\theta}^{\tau} - \nabla_{\theta} J(\theta_{\tau})\|^{2} - \frac{1}{2} \|g_{\theta}^{\tau}\|^{2} \right) \\ &\leq 3 (\varepsilon_{data} + \varepsilon_{func} + \varepsilon_{reg})^{2} + \frac{2}{(T+1)\eta_{\theta}(1-\gamma)} + \frac{3}{\alpha(T+1)} \mathbb{E}[\|g_{\theta}^{0} - \nabla_{\theta} \mathcal{L}_{0}^{D}\|^{2}] \\ &+ \frac{6\alpha\sigma^{2}}{|B|} + \frac{1}{T+1} (\frac{6\beta L^{2}}{1-\beta} + 108C_{w,Q}) d^{2} \\ &\leq 3 (\varepsilon_{data} + \varepsilon_{func} + \varepsilon_{reg})^{2} + \underbrace{(\frac{6L^{2}}{1-\beta} + \frac{108C_{w,Q}}{\alpha(2-\alpha)(1-\beta)})c}_{p_{0}} \\ &+ \underbrace{\frac{2}{T\eta_{\theta}(1-\gamma)}}_{p_{1}} + \underbrace{\frac{3}{\alpha T} \mathbb{E}[\|g_{\theta}^{0} - \nabla_{\theta} \mathcal{L}_{0}^{D}\|^{2}]}_{p_{2}} + \underbrace{\frac{6\alpha\sigma^{2}}{|B|}}_{p_{3}} + \underbrace{\frac{1}{T} (\frac{6\beta L^{2}}{1-\beta} + 108C_{w,Q})d^{2}}_{p_{4}} \end{split}$$

Next, we want to carefully choose hyper-parameters to make sure $p_0, p_2 \leq \varepsilon^2/8, p_1, p_3, p_4 \leq \varepsilon^2/4$. We consider $\beta \leq \min\{\frac{\varepsilon^2}{L^2}, \frac{(1-\gamma)^2 \varepsilon^4}{C_{\zeta,\mu}L^2}, \frac{1}{2}(1-\alpha)^2, \alpha(1-\alpha)^2\}$. Since $0 < \alpha < 1$, we have $\beta < \frac{1}{2}$.

Control p_0 For simplicity, we directly choose $\alpha = 0.9$, while α can be other constant-level values between 0 and 1. Since $\beta < 1/2$, we know

$$p_0 \le (12L^2 + 240C_{w,Q})c$$

To make sure $p_0 \leq \varepsilon^2/8$, we need

$$c \le \frac{\varepsilon^2}{12(L^2 + 20C_{w,Q})}$$

i.e.

$$\frac{8\sigma^2}{\min\{\frac{\mu_{\zeta}\eta_{\zeta}}{4},\frac{\mu_{\xi}\eta_{\xi}}{4}\}|N|}(\frac{\eta_{\zeta}}{\mu_{\zeta}}+\frac{\eta_{\xi}}{\mu_{\xi}}) \leq \frac{\varepsilon^2}{12(L^2+20C_{w,Q})}$$

Therefore,

$$|N| \geq \frac{96(L^2 + 20C_{w,Q})\sigma^2}{\min\{\frac{\mu_{\zeta}\eta_{\zeta}}{4}, \frac{\mu_{\xi}\eta_{\xi}}{4}\}\varepsilon^2} (\frac{\eta_{\zeta}}{\mu_{\zeta}} + \frac{\eta_{\xi}}{\mu_{\xi}}) = O(\varepsilon^{-2})$$

Control p_1 Since we have two constrains on η_{θ} , first we need to make sure, if $\eta_{\theta} = \frac{1}{2L_L}$

$$p_1 = \frac{4L_J}{T(1-\gamma)} \le \frac{\varepsilon^2}{4}$$

Combining 4.2, the above implies that:

$$T \ge \frac{16L_J}{(1-\gamma)\varepsilon^2} \tag{23}$$

Secondly, to make sure constraint (22) (recall that we choose $\alpha = 0.9$):

$$p_{1} = \frac{2}{T(1-\gamma)} \left(108 \left[\frac{C_{\zeta,\mu}L^{2}\beta}{18(1-\beta)} + \frac{1}{\alpha} \left(2C_{\zeta,\mu}C_{w,Q} + H^{2}C_{Q}^{2}C_{W}^{2} \right) \right] \right)^{1/2} \\ \leq \frac{2}{T(1-\gamma)} \sqrt{\frac{6C_{\zeta,\mu}L^{2}\beta}{1-\beta}} + \frac{2}{T(1-\gamma)} \sqrt{120 \left(2C_{\zeta,\mu}C_{w,Q} + H^{2}C_{Q}^{2}C_{W}^{2} \right)} \\ \leq \frac{2}{T(1-\gamma)} \sqrt{6L^{2}C_{\zeta,\mu}\frac{(1-\gamma)^{2}\varepsilon^{4}}{C_{\zeta,\mu}L^{2}}} + \frac{2}{T(1-\gamma)} \sqrt{120 \left(2C_{\zeta,\mu}C_{w,Q} + H^{2}C_{Q}^{2}C_{W}^{2} \right)} \\ = \frac{2\sqrt{6}\varepsilon^{2}}{T} + \frac{2}{T(1-\gamma)} \sqrt{120 \left(2C_{\zeta,\mu}C_{w,Q} + H^{2}C_{Q}^{2}C_{W}^{2} \right)} \right)}$$

To make sure $p_1 \leq \frac{\varepsilon^2}{4}$, we need the above two terms less than $\frac{\varepsilon^2}{8}$ at the same time, which implies

$$T \ge 16\sqrt{3}; \quad T \ge \frac{16}{(1-\gamma)\varepsilon^2} \sqrt{120 \left(2C_{\zeta,\mu}C_{w,Q} + H^2 C_Q^2 C_W^2\right)} = O(\varepsilon^{-2})$$
(24)

Control p_2 In fact, at the beginning step, $\mathbb{E}_{B_0}[g_{\theta}^0] = \nabla_{\theta} \mathcal{L}_0^D$. Therefore,

$$p_2 = \frac{\sigma^2}{|B_0|}$$

To make sure $|B_0| \geq \frac{8\sigma^2}{\varepsilon^2}$, we just set

$$|B_0| = \frac{8\sigma^2}{\varepsilon^2}.$$
(25)

Control p_3 We want $p_3 \leq \frac{\varepsilon^2}{4}$. To do that, we add the following constraint

$$\frac{|B|}{\alpha} \ge \frac{12\sigma^2}{\varepsilon^2} \tag{26}$$

which implies that $|B| \geq \frac{12\sigma^2}{\varepsilon^2}$

Control p_4 Since $\beta \leq \{1/2, \varepsilon^2/L^2\}$, we have

$$p_4 = \frac{1}{T} \left(\frac{6\beta L^2}{1-\beta} + 108C_{w,Q}\right) d^2 \le \frac{1}{T} \left(\frac{\varepsilon^2}{L^2} \frac{6L^2}{1-1/2} + 108C_{w,Q}\right) d^2 = \frac{12\varepsilon^2}{T} + 108\frac{C_{w,Q}d^2}{T} + \frac{1}{T} \frac{1}{T} \left(\frac{\varepsilon^2}{L^2} + \frac{1}{T} \frac{1}{T} + \frac{1}{T} \frac{1}{T} \frac{1}{T} \frac{1}{T} + \frac{1}{T} \frac{1}{T$$

To make sure $p_4 \leq \frac{\varepsilon^2}{4}$, we need the above two terms individually smaller than $\frac{\varepsilon^2}{8}$

$$T \ge 96; \qquad T \ge \frac{864C_{w,Q}d^2}{\varepsilon^2} \tag{27}$$

Combine (23)-(27), we need

$$T \ge \max\{96, \frac{16L_J}{\varepsilon^2}, \frac{16}{(1-\gamma)\varepsilon^2} \sqrt{120\left(2C_{\zeta,\mu}C_{w,Q} + H^2 C_Q^2 C_W^2\right)}, \frac{864C_{w,Q}d^2}{\varepsilon^2}\}$$
$$|B| \ge \frac{12\sigma^2}{\varepsilon^2}; \quad \alpha = 0.9; \quad |B_0| = \frac{8\sigma^2}{\varepsilon^2}$$

As for K in Algorithm 3, based on Theorem E.1, we can choose $K = c_{oracle} \log(\frac{1}{\beta}) = O(\log \frac{1}{\varepsilon})$ where c_{oracle} is an independent constant.

As a result, the total computation before obtaining θ_{out} should be:

$$|B_0| + |B| \cdot T + |N| \cdot K \cdot T = O(\varepsilon^{-2}) + O(\varepsilon^{-4}) + O(\varepsilon^{-2}) \cdot O(\log \frac{1}{\varepsilon}) \cdot O(\varepsilon^{-2}) = O(\varepsilon^{-4})$$

E A Concrete Example for Saddle-Point Solver Oracle

In this section, we provide an example for the oracle in Definition 4.1, which is inspired by SVRE[22].

E.1 An additional assumption

In this sub-section, we list one additional assumption for our oracle algorithm. We will illustrate the practicality of it in Appendix F.

Assumption H. Denote N as a batch data sample according to d^D whose batch size is constant |N|. We assume there exists two constants \bar{L}_{ζ} and \bar{L}_{ξ} , such that:

$$\begin{split} \mathbb{E}_{N\sim d^{D}} [\|\nabla_{\zeta}\mathcal{L}^{N}(\theta,\zeta_{1},\xi_{1})-\nabla_{\zeta}\mathcal{L}^{N}(\theta,\zeta_{2},\xi_{2})\|^{2}+\|\nabla_{\xi}\mathcal{L}^{N}(\theta,\zeta_{1},\xi_{1})-\nabla_{\xi}\mathcal{L}^{N}(\theta,\zeta_{2},\xi_{2})\|^{2}] \\ \leq \mathbb{E}_{N\sim d^{D}} \Big[-\bar{L}_{\zeta} \Big(\nabla_{\zeta}\mathcal{L}^{N}(\theta,\zeta_{1},\xi_{1})-\nabla_{\zeta}\mathcal{L}^{N}(\theta,\zeta_{2},\xi_{2})\Big)^{\top}(\zeta_{1}-\zeta_{2}) \\ &+\bar{L}_{\xi} \Big(\nabla_{\xi}\mathcal{L}^{N}(\theta,\zeta_{1},\xi_{1})-\nabla_{\xi}\mathcal{L}^{N}(\theta,\zeta_{2},\xi_{2})\Big)^{\top}(\xi_{1}-\xi_{2})\Big] \\ \mathbb{E}_{N\sim d^{D}} [\|\nabla_{\zeta}\mathcal{L}^{N}(\theta,\zeta_{1},\xi_{1})-\nabla_{\zeta}\mathcal{L}^{N}(\theta,\zeta_{2},\xi_{2})\|^{2}+\|\nabla_{\xi}\mathcal{L}^{N}(\theta,\zeta_{1},\xi_{1})-\nabla_{\xi}\mathcal{L}^{N}(\theta,\zeta_{2},\xi_{2})\|^{2}] \\ \leq \bar{L}_{\zeta}^{2} \|\zeta_{1}-\zeta_{2}\|^{2}+\bar{L}_{\xi}^{2} \|\xi_{1}-\xi_{2}\|^{2} \end{split}$$

where we use \mathcal{L}^N to denote:

$$\begin{aligned} \mathcal{L}^{N}(\theta,\zeta,\xi) &= \frac{1}{|N|} \sum_{i=1}^{|N|} (1-\gamma)Q(s^{i},a_{0}^{i})\pi(a_{0}^{i}|s^{i})\mathbb{I}[s^{i} \in S_{0}] + w(s^{i},a^{i})\Big(r^{i} + \gamma Q(s'^{i},a'^{i})\pi(a'^{i}|s'^{i}) - Q(s^{i},a^{i})\Big) \\ &+ \frac{\lambda_{Q}}{2}Q^{2}(s^{i},a^{i}) - \frac{\lambda_{w}}{2}w^{2}(s^{i},a^{i}) \end{aligned}$$

and (s^i, a^i, r^i, s'^i) are sampled from D while $a_0^i \sim \pi(\cdot | s^i), a'^i \sim \pi(\cdot | s'^i)$.

E.2 Stochastic Variance-Reduced Extragradient with Batch Data

where \mathcal{P}_{ζ} and \mathcal{P}_{ξ} are projection operator; $\nabla \mathcal{L}^{N}(\theta, \zeta, \xi)$ denotes the average gradient over samples from batch data N. Besides, we define:

$$\begin{aligned} &d^N_\zeta(\zeta_1,\xi_1,\zeta_2,\xi_2) = \nabla_\zeta \mathcal{L}^N(\theta,\zeta_1,\xi_1) - \nabla_\zeta \mathcal{L}^N(\theta,\zeta_2,\xi_2) \\ &d^N_\xi(\zeta_1,\xi_1,\zeta_2,\xi_2) = \nabla_\xi \mathcal{L}^N(\theta,\zeta_1,\xi_1) - \nabla_\xi \mathcal{L}^N(\theta,\zeta_2,\xi_2) \end{aligned}$$

Obviously,

$$\mathbb{E}[g_k^{\zeta}] = \nabla_{\zeta} \mathcal{L}^D(\theta, \zeta_k, \xi_k), \qquad \mathbb{E}[g_{k+1/2}^{\zeta}] = \nabla_{\zeta} \mathcal{L}^D(\theta, \zeta_{k+1/2}, \xi_{k+1/2})$$

where the expectation only concerns the randomness of sample when computing g. The above relationship also holds if we consider gradient w.r.t. ξ .

For this Algorithm 3, we have the following theorem: **Theorem E.1.** Under Assumption C, E, F and D, in Algorithm 3, if step size and batch size satisfy

$$\eta_{\zeta} \le \frac{1}{50 \max\{\bar{L}_{\zeta}, \mu_{\zeta}\}}, \qquad \eta_{\xi} \le \frac{1}{50 \max\{\bar{L}_{\xi}, \mu_{\xi}\}}$$

after K iterations, the algorithm will return us (ζ_K, ξ_K) :

$$\mathbb{E}[\|\zeta_K - \zeta^*\|^2 + \|\xi_K - \xi^*\|^2] \le \frac{201}{100} \left(1 - \frac{\mu\eta}{4}\right)^K \mathbb{E}[\|\zeta_0 - \zeta^*\|^2 + \|\xi_0 - \xi^*\|^2] \\ + \frac{8\sigma^2}{\min\{\frac{\mu_{\zeta}\eta_{\zeta}}{4}, \frac{\mu_{\xi}\eta_{\xi}}{4}\}|N|} \left(\frac{\eta_{\zeta}}{\mu_{\zeta}} + \frac{\eta_{\xi}}{\mu_{\xi}}\right)$$

where (ζ^*, ξ^*) is the saddle point of $\mathcal{L}^D(\theta, \zeta, \xi)$ given input θ .

We defer the proof to the next sub-section.

Algorithm 3: Stochastic Variance-Reduced Extragradient with Batch Data (SVREB)

1 Input: Stopping time K; learning rates η_{ζ}, η_{ξ} ; Initial weights ζ_0, ξ_0 ; Distribution $d^{\overline{D}}$; Batch size |N|. 2 Sample dataset $N_{\zeta}, N_{\xi} \sim d^D$ with batch size |N|3 $g_0^{\zeta} \leftarrow \nabla_{\zeta} \mathcal{L}^{N_{\zeta}}(\theta, \zeta_0, \xi_0)$

4 $g_0^{\xi} \leftarrow \nabla_{\xi} \mathcal{L}^{N_{\xi}}(\theta, \zeta_0, \xi_0)$ **5** $\zeta_1 \leftarrow \mathcal{P}_{\zeta}(\zeta_0 + \eta_{\zeta} g_0^{\zeta})$ 6 $\xi_1 \leftarrow \mathcal{P}_{\xi}(\xi_0 - \eta_{\xi} g_0^{\xi})$ 7 $m_1^{\zeta}, m_1^{\xi} \leftarrow \nabla_{\zeta} \mathcal{L}^{N_{\zeta}}(\theta, \zeta_0, \xi_0), \nabla_{\xi} \mathcal{L}^{N_{\xi}}(\theta, \zeta_0, \xi_0)$ **8** for k = 1, 2, ... K + 1 do Sample dataset $N_{\zeta}, N_{\xi} \sim d^D$ with batch size |N| $g_{k}^{\zeta} = m_{k}^{\zeta} + d_{\zeta}^{N_{\zeta}}(\zeta_{k}, \xi_{k}, \zeta_{k-1}, \xi_{k-1})$ $g_{k}^{\xi} = m_{k}^{\xi} + d_{\xi}^{N_{\xi}}(\zeta_{k}, \xi_{k}, \zeta_{k-1}, \xi_{k-1})$ 10 11 $\zeta_{k+1/2} = \mathcal{P}_{\zeta}(\zeta_k + \eta_{\zeta} g_k^{\zeta})$ 12
$$\begin{split} \xi_{k+1/2} &= \mathcal{P}_{\xi}(\xi_k - \eta_{\xi} g_k^{\xi}) \\ \text{Sample dataset } N'_{\zeta}, N'_{\xi} \sim d^D \text{ with batch size } |N| \end{split}$$
13 14 $g_{k+1/2}^{\zeta} = m_k^{\zeta} + d_{\zeta}^{N_{\zeta}'}(\zeta_{k+1/2}, \xi_{k+1/2}, \zeta_{k-1}, \xi_{k-1})$ 15 $g_{k+1/2}^{\xi} = m_k^{\xi} + d_{\xi}^{N_{\xi}'}(\zeta_{k+1/2}, \xi_{k+1/2}, \zeta_{k-1}, \xi_{k-1})$ 16 $\zeta_{k+1} = \mathcal{P}_{\zeta}(\zeta_k + \eta_{\zeta} g_{k+1/2}^{\zeta})$ 17 $\xi_{k+1} = \mathcal{P}_{\xi}(\xi_k - \eta_{\xi} g_{k+1/2}^{\xi})$ 18 // The following has been computed in step 10 and 11 19 $m_{k+1}^{\zeta}, m_{k+1}^{\xi} \leftarrow \nabla_{\zeta} \mathcal{L}^{N_{\zeta}}(\theta, \zeta_k, \xi_k), \nabla_{\xi} \mathcal{L}^{N_{\xi}}(\theta, \zeta_k, \xi_k)$ 20 $k \leftarrow k+1$ 21 22 end 23 Output: ζ_K, ξ_K

E.3 Proofs for Algorithm 3

For simplification, we will use $\omega = [\zeta, \xi] \in Z \times \Xi := \Omega$ to denote the vector concatenated by ζ and ξ . Similarly, $g_t = [-g_t^{\zeta}, g_t^{\xi}]$, and $F_N(\omega) = \mathbb{E}_{(s,a,r,s',a_0,a')\sim N}\{[-\nabla_{\zeta}\mathcal{L}^{(s,a,r,s',a_0,a')}(\theta,\zeta,\xi), \nabla_{\xi}\mathcal{L}^{(s,a,r,s',a_0,a')}(\theta,\zeta,\xi)]\}$, where N is the batch data sampled according to d^D , and $\nabla \mathcal{L}^{(s,a,r,s',a_0,a')}(\theta,\zeta,\xi)$ is the gradient computed with one sample (s, a, r, s', a_0, a') . We use $F(\omega) := \mathbb{E}_{N \sim D}[F_N(\omega)]$ to denote the gradient expected over entire dataset distribution. Besides,

$$\eta g_t = [-\eta_{\zeta} g_t^{\zeta}, \eta_{\xi} g_t^{\xi}]; \qquad \eta^2 \|\omega\|^2 = \eta_{\theta}^2 \|\zeta\|^2 + \eta_{\xi}^2 \|\xi\|^2; \qquad \mu \|\omega\|^2 = \mu_{\zeta} \|\zeta\|^2 + \mu_{\xi} \|\xi\|^2$$
$$\bar{L}^2 \|w\|^2 = \bar{L}^2_{\zeta} \|\zeta\|^2 + \bar{L}^2_{\xi} \|\xi\|^2; \qquad \eta^2 \bar{L}^2 = \eta_{\zeta}^2 \bar{L}^2_{\zeta} + \eta_{\xi}^2 \bar{L}^2_{\xi}; \qquad \eta \mu = \eta_{\zeta} \mu_{\zeta} + \eta_{\xi} \mu_{\xi}$$
The update rule for Algorithm 3 can be summarized as

The update rule for Algorithm 3 can be summarized as

Extrapolation :
$$\omega_{t+1/2} = P_{\Omega}(\omega_t - \eta g_t)$$

Update : $\omega_{t+1} = P_{\Omega}(\omega_t - \eta g_{t+1/2})$

Besides, in this section, the expectation \mathbb{E} concerns all the randomness starting from the beginning of the algorithm.

Lemma E.2 (Lemma 1 in [22]). Let $\omega \in \Omega$ and $\omega^+ := P_{\Omega}(w+u)$, then for all $w' \in \Omega$, we have $\|\omega^{+} - \omega'\|^{2} \le \|\omega - \omega'\|^{2} + 2u^{\top}(\omega^{+} - \omega') - \|\omega^{+} - \omega\|^{2}$

Lemma E.3 (Adapted from Lemma 3 in [22]). For any $w \in \Omega$, when t > 0, we have

$$\|\omega_{t+1} - \omega\|^2 \le \|\omega_t - \omega\|^2 - 2\eta g_{t+1/2}^\top (\omega_{t+1/2} - \omega) + \eta^2 \|g_t - g_{t+1/2}\| - \|\omega_{t+1/2} - \omega_t\|^2$$

and when t = 0, we have

$$\|\omega_1 - \omega\|^2 \le \|\omega_0 - \omega\|^2 - 2\eta g_0^\top (\omega_1 - \omega)$$

Proof. For t = 0, by simply applying Lemma E.2 for $(\omega, u, \omega^+, \omega') = (\omega_0, -\eta g_0^\top, \omega_1, \omega)$, we have: $\|\omega_1 - \omega\|^2 \le \|\omega_0 - \omega\|^2 - 2\eta g_0^\top(\omega_1 - \omega) - \|\omega_1 - \omega_0\|^2 \le \|\omega_0 - \omega\|^2 - 2\eta g_0^\top(\omega_1 - \omega)$ For t > 0, the proof is exactly the same as Lemma 3 in [22]

Lemma E.4 (Bound
$$||g_t - g_{t+1/2}||^2$$
). For $t > 0$, we have:

$$\mathbb{E}[||g_t - g_{t+1/2}||^2] \le 10\mathbb{E}[||F_N(w_t) - F_N(\omega^*)||^2] + 10\mathbb{E}[||F_N(\omega^*) - F_N(w_{t-1})||^2] + 5\bar{L}^2\mathbb{E}[||w_t - w_{t+1/2}||]$$

Proof. For t > 0:

$$\begin{split} & \mathbb{E}[\|g_t - g_{t+1/2}\|^2] \\ = & \mathbb{E}[\|F_N(w_t) - F_N(w_{t-1}) + m_t - F_{N'}(w_{t+1/2}) + F_{N'}(w_{t-1}) - m_t\|^2] \\ = & \mathbb{E}[\|F_N(w_t) \pm F_N(w^*) - F_N(w_{t-1}) - F_{N'}(w_{t+1/2}) \pm F_{N'}(w_t) \pm F_{N'}(w^*) + F_{N'}(w_{t-1})\|^2] \\ \leq & 5 \mathbb{E}[\|F_N(w_t) - F_N(\omega^*)\|^2] + 5 \mathbb{E}[\|F_N(\omega^*) - F_N(w_{t-1})\|^2] \\ & + 5 \mathbb{E}[\|F_{N'}(w_{t+1/2}) - F_{N'}(w_t)\|^2]] + 5 \mathbb{E}[\|F_{N'}(w_t) - F_{N'}(\omega^*)\|^2] + 5 \mathbb{E}[\|F_{N'}(\omega^*) - F_{N'}(w_{t-1})\|^2] \\ = & 10 \mathbb{E}[\|F_N(w_t) - F_N(\omega^*)\|^2] + 10 \mathbb{E}[\|F_N(\omega^*) - F_N(w_{t-1})\|^2] + 5 \mathbb{E}[\|F_{N'}(w_{t+1/2}) - F_{N'}(w_t)\|^2]] \\ & \text{where in the inequality we use the extended Young's inequality; in the last equation we use the fact that} \end{split}$$

$$\mathbb{E}_{N \sim D}[\|F_N(w_t) - F_N(w)\|^2] = \mathbb{E}_{N' \sim D}[\|F_{N'}(w_t) - F_{N'}(w)\|^2], \quad \forall w \in \Omega$$

Besides, according to Assumption H

$$\mathbb{E}[\|F_{N'}(w_{t+1/2}) - F_{N'}(w_t)\|^2]] \le \bar{L}^2 \mathbb{E}[\|w_t - w_{t+1/2}\|]$$

As a result,

$$\mathbb{E}[\|g_t - g_{t+1/2}\|^2] \le 10\mathbb{E}[\|F_N(w_t) - F_N(\omega^*)\|^2] + 10\mathbb{E}[\|F_N(\omega^*) - F_N(w_{t-1})\|^2] + 5\bar{L}^2\mathbb{E}[\|w_t - w_{t+1/2}\|]$$

Proposition E.5. Under Assumption C, for arbitrary θ , the operator $F(\omega)$ satisfying:

$$\left(F(\omega_1) - F(\omega_2)\right)^{\top} \left(\omega_1 - \omega_2\right) \ge \mu \|\omega_1 - \omega_2\|^2$$

Proof. Based on Assumption C, we have:

$$\begin{aligned} -\mathcal{L}^{D}(\theta,\zeta_{1},\xi_{2}) &\geq -\mathcal{L}^{D}(\theta,\zeta_{2},\xi_{2}) - \nabla_{\zeta}\mathcal{L}^{D}(\theta,\zeta_{1},\xi_{1})^{\top}(\zeta_{2}-\zeta_{1}) + \frac{\mu_{\zeta}}{2}\|\zeta_{2}-\zeta_{1}\|^{2} \\ -\mathcal{L}^{D}(\theta,\zeta_{2},\xi_{1}) &\geq -\mathcal{L}^{D}(\theta,\zeta_{2},\xi_{2}) - \nabla_{\zeta}\mathcal{L}^{D}(\theta,\zeta_{2},\xi_{2})^{\top}(\zeta_{1}-\zeta_{2}) + \frac{\mu_{\zeta}}{2}\|\zeta_{2}-\zeta_{1}\|^{2} \\ \mathcal{L}^{D}(\theta,\zeta_{1},\xi_{2}) &\geq \mathcal{L}^{D}(\theta,\zeta_{1},\xi_{1}) + \nabla_{\xi}\mathcal{L}^{D}(\theta,\zeta_{1},\xi_{1})^{\top}(\xi_{2}-\xi_{1}) + \frac{\mu_{\xi}}{2}\|\xi_{2}-\xi_{1}\|^{2} \\ \mathcal{L}^{D}(\theta,\zeta_{2},\xi_{1}) &\geq \mathcal{L}^{D}(\theta,\zeta_{2},\xi_{2}) + \nabla_{\xi}\mathcal{L}^{D}(\theta,\zeta_{2},\xi_{2})^{\top}(\xi_{1}-\xi_{2}) + \frac{\mu_{\xi}}{2}\|\xi_{2}-\xi_{1}\|^{2} \end{aligned}$$

Sum up and we can obtain

$$\left(F(\omega_1) - F(\omega_2) \right)^{\top} \left(\omega_1 - \omega_2 \right) := \left(F(\zeta_1, \xi_1) - F(\zeta_2, \xi_2) \right)^{\top} \left([\zeta_1, \xi_1] - [\zeta_2, \xi_2] \right)$$

= $- \left(\nabla_{\zeta} \mathcal{L}^D(\theta, \zeta_1, \xi_1) - \nabla_{\zeta} \mathcal{L}^D(\theta, \zeta_2, \xi_2) \right)^{\top} (\zeta_1 - \zeta_2) + \left(\nabla_{\xi} \mathcal{L}^D(\theta, \zeta_1, \xi_1) - \nabla_{\xi} \mathcal{L}^D(\theta, \zeta_2, \xi_2) \right)^{\top} (\xi_1 - \xi_2)$
 $\geq \mu_{\zeta} \| \zeta_2 - \zeta_1 \|^2 + \mu_{\xi} \| \xi_2 - \xi_1 \|^2 := \mu \| \omega_1 - \omega_2 \|^2$

Theorem E.1. Under Assumption C, E, F and D, in Algorithm 3, if step size and batch size satisfy

$$\eta_{\zeta} \le \frac{1}{50 \max\{\bar{L}_{\zeta}, \mu_{\zeta}\}}, \qquad \eta_{\xi} \le \frac{1}{50 \max\{\bar{L}_{\xi}, \mu_{\xi}\}}$$

after K iterations, the algorithm will return us (ζ_K, ξ_K) :

$$\mathbb{E}[\|\zeta_K - \zeta^*\|^2 + \|\xi_K - \xi^*\|^2] \le \frac{201}{100} \left(1 - \frac{\mu\eta}{4}\right)^K \mathbb{E}[\|\zeta_0 - \zeta^*\|^2 + \|\xi_0 - \xi^*\|^2] \\ + \frac{8\sigma^2}{\min\{\frac{\mu_{\zeta}\eta_{\zeta}}{4}, \frac{\mu_{\xi}\eta_{\xi}}{4}\}|N|} \left(\frac{\eta_{\zeta}}{\mu_{\zeta}} + \frac{\eta_{\xi}}{\mu_{\xi}}\right)$$

where (ζ^*, ξ^*) is the saddle point of $\mathcal{L}^D(\theta, \zeta, \xi)$ given input θ .

Proof. When t > 0, from Lemma E.3, we have

$$\|\omega_{t+1} - \omega^*\|^2 \le \|\omega_t - \omega^*\|^2 - 2\eta g_{t+1/2}^\top (\omega_{t+1/2} - \omega^*) - \|\omega_{t+1/2} - \omega_t\|^2 + \eta^2 \|g_t - g_{t+1/2}\|^2$$

Next, we use P_{t+1} to denote $\mathbb{E}[\|\omega_{t+1} - \omega^*\|^2] + \tau \mathbb{E}[\|F_N(\omega^*) - F_N(w_t)\|^2]$, where τ will be determined later, then we have

$$\begin{split} P_{t+1} &= \mathbb{E}[\|\omega_{t+1} - \omega^*\|^2] + \tau \mathbb{E}[\|F_N(\omega^*) - F_N(w_t)\|^2] \\ &\leq \mathbb{E}[\|\omega_t - \omega^*\|^2] - 2\eta \mathbb{E}[F(\omega_{t+1/2})^\top (\omega_{t+1/2} - \omega^*)] - \mathbb{E}[\|\omega_{t+1/2} - \omega_t\|^2] \\ &+ \eta^* \mathbb{E}[\|g_t - g_{t+1/2}\|^2] + \tau \mathbb{E}[\|F_N(\omega^*) - F_N(w_t)\|^2] \\ &- 2\eta \mathbb{E}[(m_t - F(\omega_{t-1}))^\top (\omega_{t+1/2} - \omega^*)] = \mathbb{E}[(F(\omega_{t+1/2}) - F(\omega_{t-1}) + m_t)^\top (\omega_{t+1/2} - \omega^*)]) \\ &\leq \mathbb{E}[\|\omega_t - \omega^*\|^2] - 2\eta \mathbb{E}[F(\omega_{t+1/2})^\top (\omega_{t+1/2} - \omega^*)] - (1 - 5\eta^2 \bar{L}^2) \mathbb{E}[\|\omega_{t+1/2} - \omega_t\|^2] \\ &+ (\tau + 10\eta^2) \mathbb{E}[\|F_N(w_t) - F_N(\omega^*)\|^2] + 10\eta^2 \mathbb{E}[\|F_N(\omega^*) - F_N(w_{t-1})\|^2] \\ &+ (\tau + 10\eta^2) \mathbb{E}[\|F_N(w_t) - F_N(\omega^*)\|^2] + 10\eta^2 \mathbb{E}[\|F_N(\omega^*) - F_N(w_{t-1})\|^2] \\ &+ 2\eta \sqrt{\mathbb{E}}[\|m_t - F(\omega_{t-1})\|^2] \mathbb{E}[\|\omega_{t+1/2} - \omega^*\|^2] \\ &\qquad (Lemma E.4 and Cauthy Inequality: \mathbb{E}[a^\top b]c] \leq \sqrt{\mathbb{E}}[\|\omega_{t+1/2} - \omega_t\|^2] \\ &+ (\tau + 10\eta^2) \mathbb{E}[\|F_N(w_t) - F_N(\omega^*)\|^2] + 10\eta^2 \mathbb{E}[\|F_N(\omega^*) - F_N(w_{t-1})\|^2] \\ &+ (\tau + 10\eta^2) \mathbb{E}[\|F_N(w_t) - F_N(\omega^*)\|^2] + 10\eta^2 \mathbb{E}[\|F_N(\omega^*) - F_N(w_{t-1})\|^2] \\ &+ (\tau + 10\eta^2) \mathbb{E}[\|W_t - W_t - F_N(\omega^*)\|^2] + 10\eta^2 \mathbb{E}[\|F_N(\omega^*) - F_N(w_{t-1})\|^2] \\ &+ (\tau + 10\eta^2) \mathbb{E}[\|F_N(w_{t+1/2})^\top (\omega_{t+1/2} - \omega^*)]^2] \\ \leq \mathbb{E}[\|\omega_t - \omega^*\|^2] - 2\eta \mathbb{E}[F(\omega_{t+1/2})^\top (\omega_{t+1/2} - \omega^*)]^2] (2\sqrt{|a^\top b|} \leq ||a||^2 + ||b||^2) \\ \leq \mathbb{E}[\|\omega_t - \omega^*\|^2] - 2\eta \mathbb{E}[F(\omega_{t+1/2})^\top (\omega_{t+1/2} - \omega^*)]^2] (2\sqrt{|a^\top b|} \leq ||a||^2 + ||b||^2) \\ \\ \leq \mathbb{E}[\|\omega_t - \omega^*\|^2] - 2\eta \mathbb{E}[F(\omega_{t+1/2})^\top (\omega_{t+1/2} - \omega^*)]^2] (1 - 25\eta^2 \bar{L}^2 - 2\tau \bar{L}^2) \mathbb{E}[\|\omega_{t+1/2} - \omega_t\|^2] \\ &+ \frac{8\sigma^2}{|N|} (\frac{\eta_c}{\mu_c} + \frac{\eta_c}{\mu_c}) + \frac{\mu\eta}{4} (\mathbb{E}[\|\omega_{t+1/2} - \omega_t\|^2 + \mathbb{E}[\|\omega_t - \omega^*\|^2]) \\ (Assumption E; Young's Inequality; \mathbb{E}[\|F_N(\omega_{t+1/2}) - F_N(\omega_t)\|^2] \leq \bar{L}^2 \mathbb{E}[\|\omega_{t+1/2} - \omega_t\|^2] \\ &+ (2\tau \bar{L} + 20\eta^2 \bar{L}) \mathbb{E}[(F_N(\omega^*) - F_N(w_{t+1/2}))^\top (\omega^* - w_{t+1/2})] + 10\eta^2 \mathbb{E}[\|F_N(\omega^*) - F_N(w_{t-1})\|^2] \\ &+ \frac{8\sigma^2}{|N|} (\frac{\eta_c}{\mu_c} + \frac{\eta_c}{\mu_c}) + \frac{\mu\eta}{4} (\mathbb{E}[\|\omega_{t+1/2} - \omega_t\|^2 + \mathbb{E}[\|\omega_t - \omega^*\|^2]) \\ (Assumption H) \\ &= \mathbb{E}[\|\omega_t - \omega^*\|^2] - (2\eta - 20\bar{L}\eta^2 - 2\tau \bar{L}\bar{L})\mathbb{E}[(F(\omega_{t+1/2}) - F(w^*))^\top (\omega_{t+1/2} - \omega^*)] \\ &- (1 - 25$$

By Prop. E.5, we have:

$$(F(w^*) - F(\omega_{t+1/2}))^{\top}(\omega^* - \omega_{t+1/2}) \ge \mu \|\omega^* - \omega_{t+1/2}\|^2 \ge \frac{\mu}{2} \|w_t - \omega^*\|^2 - \mu \|w_{t+1/2} - w_t\|^2$$
(28)

By choosing $0 < \eta_{\zeta} \leq \frac{1}{50 \max\{\bar{L}_{\zeta}, \mu_{\zeta}\}}, 0 < \eta_{\xi} \leq \frac{1}{50 \max\{\bar{L}_{\xi}, \mu_{\xi}\}}$, and $\tau = 15\eta^2$, we know $2\eta - 20\bar{L}\eta^2 - 2\tau\bar{L} = 2\eta - 50\bar{L}\eta^2 \geq 0$

As a result, we can use (28) to get:

$$\begin{aligned} P_{t+1} &\leq (1 - \mu\eta + 10\mu\eta^{2}\bar{L} + \tau\mu\bar{L} + \frac{\mu\eta}{4})\mathbb{E}[\|\omega_{t} - \omega^{*}\|^{2}] \\ &- (1 - 25\eta^{2}\bar{L}^{2} - 2\tau\bar{L}^{2} - 2\mu\eta + 20\mu\bar{L}\eta^{2} + 2\mu\tau\bar{L} - \frac{\mu\eta}{4})\mathbb{E}[\|\omega_{t+1/2} - \omega_{t}\|^{2}] \\ &+ \frac{10\eta^{2}}{\tau}\tau\mathbb{E}[\|F_{N}(\omega^{*}) - F_{N}(w_{t-1})\|^{2}] + \frac{8\sigma^{2}}{|N|}(\frac{\eta_{\zeta}}{\mu_{\zeta}} + \frac{\eta_{\xi}}{\mu_{\xi}}) \\ &\leq \underbrace{(1 - \frac{3}{4}\mu\eta + 25\mu\eta^{2}\bar{L})}_{p_{1}}\mathbb{E}[\|\omega_{t} - \omega^{*}\|^{2}] + \underbrace{(55\eta^{2}\bar{L}^{2} + \frac{9}{4}\mu\eta - 50\mu\bar{L}\eta^{2} - 1)}_{p_{2}}\mathbb{E}[\|\omega_{t+1/2} - \omega_{t}\|^{2}] \\ &+ \frac{2}{3}\tau\mathbb{E}[\|F_{N}(\omega^{*}) - F_{N}(w_{t-1})\|^{2}] + \frac{8\sigma^{2}}{|N|}(\frac{\eta_{\zeta}}{\mu_{\zeta}} + \frac{\eta_{\xi}}{\mu_{\xi}}) \end{aligned}$$

since $0 < \eta \mu \leq 1/50$ and $0 < \eta \bar{L} \leq 1/50$

$$p_1 \le 1 - \frac{3}{4}\mu\eta + \frac{25\mu\eta}{50} = 1 - \frac{\mu\eta}{4}$$
$$p_2 \le \frac{11}{500} + \frac{9}{200} - 1 \le \frac{11}{500} + \frac{9}{200} - 1 \le 0$$

As a result

$$P_{t+1} \leq (1 - \frac{\mu\eta}{4}) \mathbb{E}[||w_t - \omega^*||^2] + \frac{2}{3} \tau \mathbb{E}[||F_N(\omega^*) - F_N(w_{t-1})||^2] + \frac{8\sigma^2}{|N|} (\frac{\eta_{\zeta}}{\mu_{\zeta}} + \frac{\eta_{\xi}}{\mu_{\xi}})$$

$$\leq (1 - \min\{\frac{\mu\eta}{4}, \frac{1}{3}\}) P_t + \frac{8\sigma^2}{|N|} (\frac{\eta_{\zeta}}{\mu_{\zeta}} + \frac{\eta_{\xi}}{\mu_{\xi}})$$

$$= (1 - \frac{\mu\eta}{4}) P_t + \frac{8\sigma^2}{|N|} (\frac{\eta_{\zeta}}{\mu_{\zeta}} + \frac{\eta_{\xi}}{\mu_{\xi}})$$

$$\leq (1 - \frac{\mu\eta}{4})^t P_1 + \frac{8\sigma^2}{\min\{\frac{\mu_{\zeta}\eta_{\zeta}}{4}, \frac{\mu_{\xi}\eta_{\xi}}{4}\}|N|} (\frac{\eta_{\zeta}}{\mu_{\zeta}} + \frac{\eta_{\xi}}{\mu_{\xi}})$$

$$= (1 - \frac{\mu\eta}{4})^t (\mathbb{E}[||w_1 - \omega^*||] + \tau \mathbb{E}[||F_N(\omega^*) - F_N(\omega_0)||^2]) + \frac{8\sigma^2}{\min\{\frac{\mu_{\zeta}\eta_{\zeta}}{4}, \frac{\mu_{\xi}\eta_{\xi}}{4}\}|N|} (\frac{\eta_{\zeta}}{\mu_{\zeta}} + \frac{\eta_{\xi}}{\mu_{\xi}})$$

Next, we take a look at $\mathbb{E}[\|w_1-\omega^*\|],$ from Lemma E.3, we have:

$$\begin{split} & \mathbb{E}[\|\omega_{1}-\omega^{*}\|^{2}] \\ \leq \mathbb{E}[\|\omega_{0}-\omega^{*}\|^{2}-2\eta g_{0}^{\top}(\omega_{1}-\omega^{*})] \\ = & \mathbb{E}[\|\omega_{0}-\omega^{*}\|^{2}]+2\eta \mathbb{E}[(F(\omega_{0})-g_{0})^{\top}(\omega_{1}-\omega^{*})]+2\eta \mathbb{E}[(F(\omega^{*})-F(\omega_{0}))^{\top}(\omega_{1}-\omega^{*})] \\ \leq & \mathbb{E}[\|\omega_{0}-\omega^{*}\|^{2}]+2\eta \mathbb{E}[\|F(\omega_{0})-g_{0}\|^{2}] \mathbb{E}[\|\omega_{1}-\omega^{*}\|^{2}]+2\eta \mathbb{E}[(F(\omega^{*})-F(\omega_{0}))^{\top}(\omega_{1}-\omega^{*})] \\ \leq & \mathbb{E}[\|\omega_{0}-\omega^{*}\|^{2}]+\frac{2\eta\sigma^{2}}{|N|} \mathbb{E}[\|\omega_{1}-\omega^{*}\|]+3\eta^{2} \mathbb{E}[\|F(\omega^{*})-F(\omega_{0})\|^{2}]+\frac{1}{3} \mathbb{E}[\|\omega_{1}-\omega^{*}\|^{2}] \\ \leq & \mathbb{E}[\|\omega_{0}-\omega^{*}\|^{2}]+\frac{\mu\eta}{4} \mathbb{E}[\|\omega_{1}-\omega^{*}\|]+3\eta^{2} \mathbb{E}[\|F(\omega^{*})-F(\omega_{0})\|^{2}]+\frac{1}{3} \mathbb{E}[\|\omega_{1}-\omega^{*}\|^{2}] \\ \leq & \mathbb{E}[\|\omega_{0}-\omega^{*}\|^{2}]+\frac{1}{2} \mathbb{E}[\|\omega_{1}-\omega^{*}\|]+3\eta^{2} \mathbb{E}[\|F(\omega^{*})-F(\omega_{0})\|^{2}] \qquad (\eta\mu < 1/50) \end{split}$$

Therefore,

$$\mathbb{E}[\|\omega_1 - \omega\|^2] \le 2\mathbb{E}[\|\omega_0 - \omega^*\|^2] + 6\eta^2 \mathbb{E}[\|F(\omega^*) - F(\omega_0)\|^2]$$

Finally, using the fact that $\mathbb{E}[||F_N(\omega^*) - F_N(\omega_0)||^2] \leq \overline{L}^2 \mathbb{E}[||\omega^* - \omega_0||^2]$, we have

$$\begin{split} \mathbb{E}[\|\omega_{t+1} - \omega^*\|^2] &- \frac{8\sigma^2}{\min\{\frac{\mu_{\zeta}\eta_{\zeta}}{4}, \frac{\mu_{\xi}\eta_{\xi}}{4}\}|N|} (\frac{\eta_{\zeta}}{\mu_{\zeta}} + \frac{\eta_{\xi}}{\mu_{\xi}}) \le P_{t+1} - \frac{8\sigma^2}{\min\{\frac{\mu_{\zeta}\eta_{\zeta}}{4}, \frac{\mu_{\xi}\eta_{\xi}}{4}\}|N|} (\frac{\eta_{\zeta}}{\mu_{\zeta}} + \frac{\eta_{\xi}}{\mu_{\xi}}) \\ \le &\left(1 - \frac{\mu\eta}{4}\right)^t (2\mathbb{E}[\|\omega_0 - \omega^*\|^2] + 6\eta^2 \mathbb{E}[\|F(\omega^*) - F(\omega_0)\|^2] + \tau \mathbb{E}[\|F_N(\omega^*) - F_N(\omega_0)\|^2]) \\ = &\left(1 - \frac{\mu\eta}{4}\right)^t (2\mathbb{E}[\|\omega_0 - \omega^*\|^2] + 6\eta^2 \bar{L}^2 \mathbb{E}[\|\omega^* - \omega_0\|^2] + 15\eta^2 \bar{L}^2 \mathbb{E}[\|\omega^* - \omega_0\|^2]) \\ \le &\left(1 - \frac{\mu\eta}{4}\right)^t (2 + \frac{6}{2500} + \frac{15}{2500}) \mathbb{E}[\|\omega^* - \omega_0\|^2] \\ \le &\frac{201}{100} \left(1 - \frac{\mu\eta}{4}\right)^t \mathbb{E}[\|\omega^* - \omega_0\|^2] \end{split}$$

which finishes the proof.

F Practicality of the Assumptions in Section 2.2

First, it is common to use policy classes whose first and second order derivatives are bounded [15, 16], so the Assumption A-(1) is a reasonable one. Also, Assumption B is a common assumption in batch RL that guarantees exploratory dataset [23], and the smoothness Assumption A-(c) is frequently considered in optimization literatures.

The remaining assumptions are indeed quite strong. That said, below we show that when W and Q are the same linear class, we can satisfy these assumptions relatively easily. Indeed, Uehara et al. [4] showed that MIS-based OPE reduce to the familiar off-policy LSTD algorithms with linear classes [24, 25], and we show that Assumptions A-(b), C, D, E, F, G can be satisfied in this case if we simply assume Assumption I, which is standard in the off-policy LSTD literature.

Definition F.1 (Linear function classes). We have a feature class $\{\phi(s, a) \in \mathbb{R}^{u \times 1} | \forall s, a \in S \times A\}$ subject to $\|\phi(s, a)\| = 1$, and two parameter spaces $Z, \Xi \in \mathbb{R}^{u \times 1}$. The approximated value function Q_{ξ} and density ratio w_{ζ} are represented by

$$w(\cdot, \cdot) = \boldsymbol{\phi}(\cdot, \cdot)^{\top} \boldsymbol{\zeta}, \quad Q(\cdot, \cdot) = \boldsymbol{\phi}(\cdot, \cdot)^{\top} \boldsymbol{\xi}$$

Remark F.2. Since $\|\phi(\cdot, \cdot)\| \leq 1$, the matrix $\mathbb{E}_{s,a\sim d^D}[\phi(s,a)\phi(s,a)^{\top}]$ is semi-positive definite and its largest eigenvalue is less than 1.

Assumption I. There exists a positive constant σ_{\min} that, the matrix $\mathbb{E}_{s,a\sim d^D}[\phi(s,a)\phi(s,a)^{\top}]$ is full-rank, and all its eigenvalues are no less than σ_{\min} ; besides, the matrix $\mathbb{E}_{s,a,s',a'\sim d^D}[\phi(s,a)\phi(s,a)^{\top} - \gamma\phi(s,a)\phi(s',a')]$ is invertible, and its minimal sigular value is no less than σ_{\min} .

Remark F.3. In Assumption I, we only add requirement on the smallest singular value of \mathbf{M} and do not care about whether all its eigenvalues are positive or not.

For simplicity, we choose $\lambda_w = \lambda_Q = \lambda > 0$. We use $\Phi \in \mathbb{R}^{|S||A| \times u}$ to denote the matrix concatenated by all features, use **K** to denote $\Phi^{\top} \Lambda^D \Phi$ and use **M** to denote $\Phi^{\top} \Lambda^D (I - \gamma \mathbf{P}_D^{\pi}) \Phi$, where Λ^D is a diagonal matrix whose diagonal elements are $d^D(\cdot, \cdot)$, and \mathbf{P}_D^{π} is the empirical transition matrix induced from dataset distribution. Notes that if we never see some *s* in dataset, then the corresponding element in Λ^D should be 0, and we do not need to worry about the corresponding row in \mathbf{P}_D^{π} . By choosing linear function classes, we can rewrite \mathcal{L}^D to:

$$\mathcal{L}^{D}(\pi,\zeta,\xi) = (1-\gamma)\mathbb{E}_{s_{0}}[Q(s_{0},\pi)] + \mathbb{E}_{w}[r+\gamma Q(s',\pi) - Q(s,a)] + \frac{\lambda}{2}\mathbb{E}_{d^{D}}[Q^{2}(s,a)] - \frac{\lambda}{2}\mathbb{E}_{d^{D}}[w^{2}(s,a)]$$
$$= (1-\gamma)\nu_{D}^{\pi}\Phi\xi + \zeta^{\top}\Phi^{\top}\Lambda^{D}(R - (I-\gamma\mathbf{P}_{D}^{\pi})\Phi\xi) + \frac{\lambda}{2}\xi^{\top}\mathbf{K}\xi - \frac{\lambda}{2}\zeta^{\top}\mathbf{K}\zeta$$
$$= (1-\gamma)\nu_{D}^{\pi}\Phi\xi + \zeta^{\top}\Phi^{\top}\Lambda^{D}R - \zeta^{\top}\mathbf{M}\xi + \frac{\lambda}{2}\xi^{\top}\mathbf{K}\xi - \frac{\lambda}{2}\zeta^{\top}\mathbf{K}\zeta$$

Since \mathcal{L}^D is quadratic, under Assumption I, matrix **K** is full-rank with minimal eigenvalue larger than σ_{\min} and maximal eigenvalue smaller than 1, then $\mathcal{L}^D(\pi, \zeta, \xi)$ is $\lambda \sigma_{\min}$ -strongly-concave- $\lambda \sigma_{\min}$ -strongly-convex, and λ smooth. Combining bounded second order derivatives of $\log \pi$, \mathcal{L} is also smooth w.r.t. θ . Therefore, we know Assumption C holds.

Next, we try to give a bound for the norm of the saddle point of $\mathcal{L}^D(\pi, w_{\zeta}, Q_{\xi})$ denotes as (ζ^*, ξ^*) , to testify the other assumptions. By taking derivatives w.r.t. ξ , we have:

$$\xi = \frac{1}{\lambda} \mathbf{K}^{-1} \Big(\mathbf{M}^{\top} \zeta - (1 - \gamma) \boldsymbol{\Phi}^{\top} (\boldsymbol{\nu}_{D}^{\pi})^{\top} \Big)$$

Plug it into \mathcal{L}^D :

$$-\frac{\lambda}{2}\zeta^{\top}\mathbf{K}\zeta - \frac{1}{2\lambda}\left(\mathbf{M}^{\top}\zeta - (1-\gamma)\boldsymbol{\Phi}^{\top}(\nu_{D}^{\pi})^{\top}\right)^{\top}\mathbf{K}^{-1}\left(\mathbf{M}^{\top}\zeta - (1-\gamma)\boldsymbol{\Phi}^{\top}(\nu_{D}^{\pi})^{\top}\right) + \zeta^{\top}\boldsymbol{\Phi}^{\top}\boldsymbol{\Lambda}^{D}R$$

Taking the derivative of ζ , we have:

Taking the derivative of ζ , we have:

$$\zeta^* = \left(\lambda^2 \mathbf{K} + \mathbf{M} \mathbf{K}^{-1} \mathbf{M}^{\top}\right)^{-1} \left(-(1-\gamma) \mathbf{M} \mathbf{K}^{-1} \mathbf{\Phi}^{\top} (\nu_D^{\pi})^{\top} + \lambda \mathbf{\Phi}^{\top} \mathbf{\Lambda}^D R\right)$$

and therefore,

$$\begin{split} \xi^* &= \frac{1}{\lambda} \mathbf{K}^{-1} \Big(\mathbf{M}^\top \zeta^* - (1 - \gamma) \mathbf{\Phi}^\top (\nu_D^\pi)^\top \Big) \\ &= \frac{1}{\lambda} \mathbf{K}^{-1} \mathbf{M}^\top \Big(\lambda^2 \mathbf{K} + \mathbf{M} \mathbf{K}^{-1} \mathbf{M}^\top \Big)^{-1} \cdot \Big(- (1 - \gamma) \mathbf{M} \mathbf{K}^{-1} \mathbf{\Phi}^\top (\nu_D^\pi)^\top + \lambda \mathbf{\Phi}^\top \mathbf{\Lambda}^D R \Big) \\ &+ (1 - \gamma) \frac{1}{\lambda} \mathbf{K}^{-1} \mathbf{\Phi}^\top (\nu_D^\pi)^\top \\ &= (1 - \gamma) \lambda \Big(\lambda^2 \mathbf{K} + \mathbf{M}^\top \mathbf{K}^{-1} \mathbf{M} \Big)^{-1} \cdot \mathbf{\Phi}^\top (\nu_D^\pi)^\top + \mathbf{K}^{-1} \mathbf{M}^\top \Big(\lambda^2 \mathbf{K} + \mathbf{M} \mathbf{K}^{-1} \mathbf{M}^\top \Big)^{-1} \mathbf{\Phi}^\top \mathbf{\Lambda}^D R \end{split}$$

where in the last step, we use the inverse matrix lemma:

$$(\lambda^{2}\mathbf{K} + \mathbf{M}^{\top}\mathbf{K}^{-1}\mathbf{M})^{-1} = \frac{1}{\lambda^{2}}\mathbf{K}^{-1} - \frac{1}{\lambda^{2}}\mathbf{K}^{-1}\mathbf{M}^{\top}(\lambda^{2}\mathbf{K} + \mathbf{M}\mathbf{K}^{-1}\mathbf{M}^{\top})\mathbf{M}\mathbf{K}^{-1}$$

Because $\|\phi(\cdot, \cdot)\| \leq 1$, it's easy to prove that, for arbitrary vector $x \in \mathbb{R}^d$,

$$\max\{\|\mathbf{M}x\|, \|\mathbf{M}^{\top}x\|\} \le (1+\gamma)\|x\|$$

Therefore,

$$\begin{split} \|\zeta^*\| &\leq (1-\gamma) \| \left(\lambda^2 \mathbf{K} + \mathbf{M} \mathbf{K}^{-1} \mathbf{M}^\top \right)^{-1} \mathbf{M} \mathbf{K}^{-1} \| \cdot \| \mathbf{\Phi}^\top (\nu_D^\pi)^\top \| + \| \left(\lambda^2 \mathbf{K} + \mathbf{M} \mathbf{K}^{-1} \mathbf{M}^\top \right)^{-1} \| \cdot \| \lambda \mathbf{\Phi}^\top \mathbf{\Lambda}^D R \| \\ &\leq (1-\gamma) \| \left(\mathbf{M}^\top \right)^{-1} \| + \lambda \| \left(\lambda^2 \mathbf{K} + \mathbf{M} \mathbf{K}^{-1} \mathbf{M}^\top \right)^{-1} \| \\ &\leq (1-\gamma) \| \left(\mathbf{M}^\top \right)^{-1} \| + \lambda \| \left(\lambda^2 \mathbf{K} + \mathbf{M} \mathbf{K}^{-1} \mathbf{M}^\top \right)^{-1} \| \\ &\leq \frac{1-\gamma}{\sigma_{\min}} + \frac{\lambda}{\lambda^2 \sigma_{\min} + \sigma_{\min}^2} \coloneqq D_{\zeta} \\ \| \xi^* \| \leq (1-\gamma) \lambda \| \left(\lambda^2 \mathbf{K} + \mathbf{M}^\top \mathbf{K}^{-1} \mathbf{M} \right)^{-1} \| \cdot \| \mathbf{\Phi}^\top (\nu_D^\pi)^\top \| \\ &+ \| \mathbf{K}^{-1} \mathbf{M}^\top \left(\lambda^2 \mathbf{K} + \mathbf{M} \mathbf{K}^{-1} \mathbf{M}^\top \right)^{-1} \mathbf{M} \mathbf{K}^{-1} \mathbf{K} \mathbf{M}^{-1} \| \| \mathbf{\Phi}^\top \mathbf{\Lambda}^D R \| \\ \leq (1-\gamma) \lambda \| \left(\lambda^2 \mathbf{K} + \mathbf{M}^\top \mathbf{K}^{-1} \mathbf{M} \right)^{-1} \| + \| (\mathbf{K}^{-1} - (\mathbf{K} + \frac{1}{\lambda^2} \mathbf{M}^\top \mathbf{K}^{-1} \mathbf{M})^{-1} \| \\ \leq (1-\gamma) \lambda \| \left(\lambda^2 \mathbf{K} + \mathbf{M}^\top \mathbf{K}^{-1} \mathbf{M} \right)^{-1} \| + \| \mathbf{K}^{-1} \mathbf{K} \mathbf{M}^{-1} \| \\ \leq \frac{(1-\gamma) \lambda}{\lambda^2 \sigma_{\min} + \sigma_{\min}^2}} + \frac{1}{\sigma_{\min}} \coloneqq D_{\xi} \end{split}$$

By choosing $Z = \{\zeta | \|\zeta\| \le D_{\zeta} + 1\}$ and $\Xi = \{\xi | \|\xi\| \le D_{\xi} + 1\}$, Assumptions D and F, G can be satisfied when $d = 2 \max\{D_{\zeta}, D_{\xi}\} + 2$. Moreover,

$$w_{\zeta}(s,a) = \phi(s,a)^{\top} \zeta \leq \|\phi(s,a)\| \|\zeta\| \leq D_{\zeta}$$

$$Q_{\xi}(s,a) = \phi(s,a)^{\top} \xi \leq \|\phi(x,a)\| \|\xi\| \leq D_{\xi}$$

$$\|w_{\zeta_{1}}(s,a) - w_{\zeta_{2}}(s,a)\| \leq \|\phi(s,a)\| \|\zeta_{1} - \zeta_{2}\| \leq \|\zeta_{1} - \zeta_{2}\|$$

$$\|Q_{\xi_{1}}(s,a) - Q_{\xi_{2}}(s,a)\| \leq \|\phi(s,a)\| \|\xi_{1} - \xi_{2}\| \leq \|\xi_{1} - \xi_{2}\|$$

which means Assumption A-(b) is satisfied by setting $C_W = D_\zeta$, $C_Q = D_\xi$ and $L_W = L_Q = 1$. Besides, D_ζ and D_ξ are finite also implies that σ in Assumption E is finite.

Finally, we evaluate Assumption H. For simplicity, we use \mathbf{K}_N to denote matrix $\mathbb{E}_B[\phi(s, a)\phi(s, a)^\top]$ and use \mathbf{M}_N to denote $\mathbb{E}_B[\phi(s, a)\phi(s, a)^\top - \phi(s, a)\phi(s', \pi)^\top]$

$$\nabla_{\zeta} \mathcal{L}^{N}(\theta, \zeta_{1}, \xi_{1}) - \nabla_{\zeta} \mathcal{L}^{N}(\theta, \zeta_{2}, \xi_{2}) = -\lambda \mathbf{K}_{N}(\zeta_{1} - \zeta_{2}) - \mathbf{M}_{N}(\xi_{1} - \xi_{2})$$
$$\nabla_{\xi} \mathcal{L}^{N}(\theta, \zeta_{1}, \xi_{1}) - \nabla_{\xi} \mathcal{L}^{N}(\theta, \zeta_{2}, \xi_{2}) = \lambda \mathbf{K}_{N}(\xi_{1} - \xi_{2}) - \mathbf{M}_{N}^{\top}(\zeta_{1} - \zeta_{2})$$

Therefore,

$$\begin{split} & \mathbb{E}_{N \sim d^{D}} [\| \nabla_{\zeta} \mathcal{L}^{N}(\theta, \zeta_{1}, \xi_{1}) - \nabla_{\zeta} \mathcal{L}^{N}(\theta, \zeta_{2}, \xi_{2}) \|^{2} + \| \nabla_{\xi} \mathcal{L}^{N}(\theta, \zeta_{1}, \xi_{2}) - \nabla_{\xi} \mathcal{L}^{N}(\theta, \zeta_{2}, \xi_{2}) \|^{2}] \\ & \leq 2 \mathbb{E}_{N \sim d^{D}} [(\zeta_{1} - \zeta_{2})^{\top} (\lambda^{2} \mathbf{K}_{N}^{\top} \mathbf{K}_{N} + \mathbf{M}_{N}^{\top} \mathbf{M}_{N}) (\zeta_{1} - \zeta_{2})] \\ & + 2 \mathbb{E}_{N \sim d^{D}} [(\xi_{1} - \xi_{2})^{\top} (\lambda^{2} \mathbf{K}_{N}^{\top} \mathbf{K}_{N} + \mathbf{M}_{N}^{\top} \mathbf{M}_{N}) (\xi_{1} - \xi_{2})] \\ & \leq 2 \mathbb{E}_{N \sim d^{D}} [(\zeta_{1} - \zeta_{2})^{\top} (\lambda^{2} \mathbf{K}_{N}^{2} + (1 + \gamma)^{2} I) (\zeta_{1} - \zeta_{2})] \\ & + 2 \mathbb{E}_{N \sim d^{D}} [(\xi_{1} - \xi_{2})^{\top} (\lambda^{2} \mathbf{K}_{N}^{2} + (1 + \gamma)^{2} I) (\xi_{1} - \xi_{2})] \\ & \leq 2 \mathbb{E}_{N \sim d^{D}} [(\zeta_{1} - \zeta_{2})^{\top} (\lambda^{2} \mathbf{K}_{N} + (1 + \gamma)^{2} I) (\zeta_{1} - \zeta_{2})] \\ & + 2 \mathbb{E}_{N \sim d^{D}} [(\xi_{1} - \xi_{2})^{\top} (\lambda^{2} \mathbf{K}_{N} + (1 + \gamma)^{2} I) (\xi_{1} - \xi_{2})] \\ & = (\zeta_{1} - \zeta_{2})^{\top} (2\lambda^{2} \mathbf{K} + 2(1 + \gamma)^{2} I) (\zeta_{1} - \zeta_{2}) + (\xi_{1} - \xi_{2})^{\top} (2\lambda^{2} \mathbf{K} + 2(1 + \gamma)^{2} I) (\xi_{1} - \xi_{2}) \end{split}$$

In the first inequality, we use Young's inequality; in the second one, we use the fact that the largest singular value of \mathbf{M}_N is less than $(1 + \gamma)$; the third one is because all eigenvalues of \mathbf{K}_N locate in [0, 1], and we should have $I \succ \mathbf{K}_N \succ \mathbf{K}_N^2$. Notice that,

$$\mathbb{E}_{N\sim d^{D}} \left[-\left(\nabla_{\zeta} \mathcal{L}^{N}(\theta, \zeta_{1}, \xi_{1}) - \nabla_{\zeta} \mathcal{L}^{N}(\theta, \zeta_{2}, \xi_{2}) \right)^{\top} (\zeta_{1} - \zeta_{2}) + \left(\nabla_{\xi} \mathcal{L}^{N}(\theta, \zeta_{1}, \xi_{1}) - \nabla_{\xi} \mathcal{L}^{N}(\theta, \zeta_{2}, \xi_{2}) \right)^{\top} (\xi_{1} - \xi_{2}) \right]$$

$$\approx \lambda (\zeta_{1} - \zeta_{2})^{\top} \mathbf{K} (\zeta_{1} - \zeta_{2}) + \lambda (\xi_{1} - \xi_{2})^{\top} \mathbf{K} (\xi_{1} - \xi_{2})$$

Therefore,

=

$$\begin{aligned} &(\zeta_1 - \zeta_2)^\top (2\lambda^2 \mathbf{K} + 2(1+\gamma)^2 I)(\zeta_1 - \zeta_2) + (\xi_1 - \xi_2)^\top (2\lambda^2 \mathbf{K} + 2(1+\gamma)^2 I)(\xi_1 - \xi_2) \\ &\leq (2\lambda + \frac{2(1+\gamma)^2}{\sigma_{\min}\lambda}) \Big(\lambda(\zeta_1 - \zeta_2)^\top \mathbf{K}(\zeta_1 - \zeta_2) + \lambda(\xi_1 - \xi_2)^\top \mathbf{K}(\xi_1 - \xi_2)\Big) \end{aligned}$$

Moreover,

$$(\zeta_1 - \zeta_2)^\top (2\lambda^2 \mathbf{K} + 2(1+\gamma)^2 I)(\zeta_1 - \zeta_2) + (\xi_1 - \xi_2)^\top (2\lambda^2 \mathbf{K} + 2(1+\gamma)^2 I)(\xi_1 - \xi_2)$$

$$\leq (2\lambda^2 + 2(1+\gamma)^2) \Big((\zeta_1 - \zeta_2)^\top (\zeta_1 - \zeta_2) + (\xi_1 - \xi_2)^\top (\xi_1 - \xi_2) \Big)$$

As a result, Assumption H holds with $\bar{L}_{\zeta} = \bar{L}_{\xi} = \max\{2\lambda + \frac{2(1+\gamma)^2}{\sigma_{\min}\lambda}, \sqrt{2\lambda^2 + 2(1+\gamma)^2}\}.$