

## A Pseudocode of TD Learning

In this section, we present the pseudocode of TD learning in Algorithm 1, which is introduced in §3.

---

### Algorithm 1 Temporal-Difference Learning with Two-Layer Neural Network for Policy Evaluation

---

**Initialization:**  $\theta_i(0) \stackrel{\text{i.i.d.}}{\sim} \rho_0$  ( $i \in [m]$ ), number of iterations  $K = \lfloor T/\epsilon \rfloor$ , and policy  $\pi$  of interest.  
**for**  $k = 0, \dots, K - 1$  **do**  
    Sample the state-action pair  $(s, a)$  from the stationary distribution  $\mathcal{D}$  of  $\pi$ , receive the reward  $r$ , and obtain the subsequent state-action pair  $(s', a')$ .  
    Calculate the Bellman residual  $\delta = \widehat{Q}(x; \theta^{(m)}(k)) - r - \gamma \cdot \widehat{Q}(x'; \theta^{(m)}(k))$ , where  $x = (s, a)$  and  $x' = (s', a')$ .  
    Perform the TD update  $\theta_i(k+1) \leftarrow \theta_i(k) - \eta \epsilon \cdot \alpha \cdot \delta \cdot \nabla_{\theta} \sigma(x; \theta_i(k))$  ( $i \in [m]$ ).  
**end for**  
**Output:**  $\{\theta^{(m)}(k)\}_{k=0}^{K-1}$

---

## B Q-Learning and Policy Improvement

In this section, we extend our analysis of TD to Q-learning and soft Q-learning for policy improvement. In §B.1, we introduce Q-learning and its mean-field limit. In §B.2, we establish the global optimality and convergence of Q-learning. In §B.3, we further extend our analysis to soft Q-learning, which is equivalent to policy gradient.

### B.1 Q-Learning

Q-learning aims to solve the following projected Bellman optimality equation,

$$Q = \Pi_{\mathcal{F}} \mathcal{T}^* Q. \quad (\text{B.1})$$

Here  $\mathcal{T}^*$  is the Bellman optimality operator, which is defined as follows,

$$\mathcal{T}^* Q(s, a) = \mathbb{E} \left[ r + \gamma \cdot \max_{\underline{a} \in \mathcal{A}} Q(s', \underline{a}) \mid r \sim R(\cdot \mid s, a), s' \sim P(\cdot \mid s, a) \right].$$

When  $\Pi_{\mathcal{F}}$  is the identity mapping, the fixed point solution to (B.1) is the Q-function  $Q^{\pi^*}$  of the optimal policy  $\pi^*$ , which maximizes the expected total reward  $J(\pi)$  defined in (2.1) [65]. We consider the parameterization of the Q-function in (3.1) and update the parameter  $\theta^{(m)}$  as follows,

$$\begin{aligned} & \theta_i(k+1) \\ &= \theta_i(k) - \eta \epsilon \cdot \alpha \cdot \left( \widehat{Q}(s_k, a_k; \theta^{(m)}(k)) - r_k - \gamma \cdot \max_{\underline{a} \in \mathcal{A}} \widehat{Q}(s'_k, \underline{a}; \theta^{(m)}(k)) \right) \cdot \nabla_{\theta} \sigma(s_k, a_k; \theta_i(k)), \end{aligned} \quad (\text{B.2})$$

where  $i \in [m]$ ,  $(s_k, a_k)$  is sampled from the stationary distribution  $\mathcal{D}_{\text{E}} \in \mathcal{P}(\mathcal{S} \times \mathcal{A})$  of an exploration policy  $\pi_{\text{E}}$ ,  $r_k \sim R(\cdot \mid s_k, a_k)$  is the reward, and  $s'_k \sim P(\cdot \mid s_k, a_k)$  is the subsequent state. For notational simplicity, we denote by  $\widetilde{\mathcal{D}}_{\text{E}} \in \mathcal{P}(\mathcal{S} \times \mathcal{A} \times \mathbb{R} \times \mathcal{S})$  the distribution of  $(s_k, a_k, r_k, s'_k)$ . For an initial distribution  $\nu_0 \in \mathcal{P}(\mathbb{R}^D)$ , we initialize  $\{\theta_i\}_{i=1}^m$  as  $\theta_i \stackrel{\text{i.i.d.}}{\sim} \rho_0$  ( $i \in [m]$ ). See Algorithm 2 for a detailed description.

**Mean-Field Limit.** Corresponding to  $\epsilon \rightarrow 0^+$  and  $m \rightarrow \infty$ , the mean-field limit of the Q-learning dynamics in (B.2) is characterized by the following PDE with  $\nu_0$  as the initial distribution,

$$\partial_t \nu_t = -\eta \cdot \text{div}(\nu_t \cdot h(\cdot; \nu_t)). \quad (\text{B.3})$$

Here  $h(\cdot; \nu_t) : \mathbb{R}^D \rightarrow \mathbb{R}^D$  is a vector field, which is defined as follows,

$$h(\theta; \nu) = -\alpha \cdot \mathbb{E}_{(s, a, r, s') \sim \widetilde{\mathcal{D}}_{\text{E}}} \left[ \left( Q(s, a; \nu) - r - \gamma \cdot \max_{\underline{a} \in \mathcal{A}} Q(s', \underline{a}; \nu) \right) \cdot \nabla_{\theta} \sigma(s, a; \theta) \right]. \quad (\text{B.4})$$

In parallel to Proposition 3.1, the empirical distribution  $\widehat{\nu}_k^{(m)} = m^{-1} \cdot \sum_{i=1}^m \delta_{\theta_i(k)}$  weakly converges to  $\nu_{k\epsilon}$  as  $\epsilon \rightarrow 0^+$  and  $m \rightarrow \infty$ .

---

**Algorithm 2** Q-Learning with Two-Layer Neural Network for Policy Improvement

---

**Initialization.**  $\theta_i(0) \stackrel{\text{i.i.d.}}{\sim} \nu_0$  ( $i \in [m]$ ), number of iterations  $K = \lfloor T/\epsilon \rfloor$ , and exploration policy  $\pi_E$ .

**for**  $k = 0, \dots, K - 1$  **do**

Sample the state-action pair  $(s, a)$  from the stationary distribution  $\mathcal{D}_E$  of  $\pi_E$ , receive the reward  $r$ , and obtain the subsequent state  $s'$ .

Calculate the Bellman residual  $\delta = \widehat{Q}(x; \theta^{(m)}(k)) - r - \gamma \cdot \widehat{Q}(x'; \theta^{(m)}(k))$ , where  $x = (s, a)$  and  $x' = (s', \arg\max_{a \in \mathcal{A}} \widehat{Q}(s', a; \theta^{(m)}(k)))$ .

Perform the Q-learning update  $\theta_i(k+1) \leftarrow \theta_i(k) - \eta \epsilon \cdot \alpha \cdot \delta \cdot \nabla_{\theta} \sigma(x; \theta_i(k))$  ( $i \in [m]$ ).

**end for**

**Output:**  $\{\theta^{(m)}(k)\}_{k=0}^{K-1}$

---

## B.2 Global Optimality and Convergence of Q-Learning

The  $\max$  operator in the Bellman optimality operator  $\mathcal{T}^*$  makes the analysis of Q-learning more challenging than that of TD. Correspondingly, we lay out an extra regularity condition on the exploration policy  $\pi_E$ . Recall that the function class  $\mathcal{F}$  is defined in (4.3).

**Assumption B.1.** We assume for an absolute constant  $\kappa > 0$  and any  $Q^1, Q^2 \in \mathcal{F}$  that

$$\mathbb{E}_{(s,a) \sim \mathcal{D}_E} \left[ (Q^1(s, a) - Q^2(s, a))^2 \right] \geq (\gamma + \kappa)^2 \cdot \mathbb{E}_{(s,a) \sim \mathcal{D}_E} \left[ \left( \max_{\underline{a} \in \mathcal{A}} Q^1(s, \underline{a}) - \max_{\underline{a} \in \mathcal{A}} Q^2(s, \underline{a}) \right)^2 \right].$$

Although Assumption B.1 is strong, we are not aware of any weaker regularity condition in the literature, even in the linear setting [25, 55, 78] and the NTK regime [21]. Let the initial distribution  $\nu_0$  be the standard Gaussian distribution  $N(0, I_D)$ . In parallel to Theorem 4.3, we establish the following theorem, which characterizes the global optimality and convergence of Q-learning. Recall that we write  $\mathcal{X} = \mathcal{S} \times \mathcal{A}$  and  $x = (s, a) \in \mathcal{X}$ . Also,  $\nu_t$  is the PDE solution in (B.3), while  $\theta^{(m)}(k)$  is the Q-learning dynamics in (B.2).

**Theorem B.2.** There exists a unique fixed point solution to the projected Bellman optimality equation  $Q = \Pi_{\mathcal{F}} \mathcal{T}^* Q$ , which takes the form of  $Q^\dagger(x) = \int \sigma(x; \theta) d\bar{\nu}(\theta)$ . We assume that  $D_{\chi^2}(\bar{\nu} \parallel \nu_0) < \infty$  and  $\bar{\nu}(\theta) > 0$  for any  $\theta \in \mathbb{R}^D$ . Under Assumptions 4.1, 4.2, and B.1, it holds for  $\eta = \alpha^{-2}$  that

$$\inf_{t \in [0, T]} \mathbb{E}_{x \sim \mathcal{D}_E} \left[ (Q(x; \nu_t) - Q^\dagger(x))^2 \right] \leq \frac{(\kappa + \gamma) \cdot D_{\chi^2}(\bar{\nu} \parallel \nu_0)}{2\kappa \cdot T} + \frac{(\kappa + \gamma) \cdot C_*}{\kappa \cdot \alpha}, \quad (\text{B.5})$$

where  $C_* > 0$  is a constant depending on  $D_{\chi^2}(\bar{\nu} \parallel \nu_0)$ ,  $B_1$ ,  $B_2$ , and  $B_r$ . Moreover, it holds with probability at least  $1 - \delta$  that

$$\begin{aligned} & \min_{\substack{k \leq T/\epsilon \\ (k \in \mathbb{N})}} \mathbb{E}_{x \sim \mathcal{D}_E} \left[ \left( \widehat{Q}(x; \theta^{(m)}(k)) - Q^\dagger(x) \right)^2 \right] \\ & \leq \frac{(\kappa + \gamma) \cdot D_{\chi^2}(\bar{\nu} \parallel \nu_0)}{2\kappa \cdot T} + \frac{(\kappa + \gamma) \cdot C_*}{\kappa \cdot \alpha} + \Delta(\epsilon, m, \delta, T), \end{aligned} \quad (\text{B.6})$$

where  $\Delta(\epsilon, m, \delta, T) > 0$  is an error term such that

$$\lim_{m \rightarrow \infty} \lim_{\epsilon \rightarrow 0^+} \Delta(\epsilon, m, \delta, T) = 0.$$

*Proof.* See §B.4 for a detailed proof. □

Theorem B.2 proves that the optimality gap  $\mathbb{E}_{x \sim \mathcal{D}_E} [(Q(x; \nu_t) - Q^\dagger(x))^2]$  decays to zero at a sublinear rate up to the error of  $O(\alpha^{-1})$ , where  $\alpha > 0$  is the scaling parameter in (3.1). In parallel to Theorem 4.3, varying  $\alpha$  leads to a tradeoff between such an error of  $O(\alpha^{-1})$  and the deviation of  $\nu_t$  from  $\nu_0$ . Moreover, based on the counterparts of Proposition 3.1 and Lemma D.6, Theorem B.2 gives the global optimality and convergence of the Q-learning dynamics  $\theta^{(m)}(k)$  in (B.2), which is in parallel to Corollary 4.4.

### B.3 Soft Q-Learning

In this section, we generalize Theorem B.2 to soft Q-learning. To introduce soft Q-learning, we first define the soft Bellman optimality operator as follows,

$$\mathcal{T}_\beta Q(s, a) = \mathbb{E} \left[ r + \gamma \cdot \operatorname{softmax}_{\underline{a} \in \mathcal{A}}^\beta Q(s', \underline{a}) \mid r \sim R(\cdot \mid s, a), s' \sim P(\cdot \mid s, a) \right],$$

where the softmax operator is defined as follows,

$$\operatorname{softmax}_{\underline{a} \in \mathcal{A}}^\beta Q(s, \underline{a}) = \beta \cdot \log \mathbb{E}_{\underline{a} \sim \bar{\pi}(\cdot \mid s)} \left[ \exp(\beta^{-1} \cdot Q(s, \underline{a})) \right].$$

Here  $\bar{\pi}(\cdot \mid s)$  is the uniform policy. Soft Q-learning aims to find the fixed point solution to the projected soft Bellman optimality equation  $Q = \Pi_{\mathcal{F}} \mathcal{T}_\beta Q$ . In parallel to the Q-learning dynamics in (B.2), we consider the following soft Q-learning dynamics,

$$\begin{aligned} \theta_i(k+1) &= \theta_i(k) - \eta \epsilon \cdot \alpha \cdot \left( \widehat{Q}(s_k, a_k; \theta^{(m)}(k)) - r_k - \gamma \cdot \operatorname{softmax}_{\underline{a} \in \mathcal{A}}^\beta \widehat{Q}(s'_k, \underline{a}; \theta^{(m)}(k)) \right) \cdot \nabla_{\theta} \sigma(s_k, a_k; \theta_i(k)), \end{aligned} \quad (\text{B.7})$$

whose mean-field limit is characterized by the following PDE,

$$\partial_t \nu_t = -\eta \cdot \operatorname{div}(\nu_t \cdot h(\cdot; \nu_t)). \quad (\text{B.8})$$

In parallel to (B.4),  $h(\cdot; \nu_t) : \mathbb{R}^D \rightarrow \mathbb{R}^D$  is a vector field, which is defined as follows,

$$h(\theta; \nu) = -\alpha \cdot \mathbb{E}_{(s, a, r, s') \sim \widehat{\mathcal{D}}_E} \left[ (Q(s, a; \nu) - r - \gamma \cdot \operatorname{softmax}_{\underline{a} \in \mathcal{A}}^\beta Q(s', \underline{a}; \nu)) \cdot \nabla_{\theta} \sigma(s, a; \theta) \right].$$

In parallel to Assumption B.1, we lay out the following regularity condition.

**Assumption B.3.** We assume for an absolute constant  $\kappa > 0$  and any  $\nu^1, \nu^2 \in \mathcal{P}(\mathbb{R}^D)$  that

$$\begin{aligned} &\mathbb{E}_{(s, a) \sim \mathcal{D}_E} \left[ (Q(s, a; \nu^1) - Q(s, a; \nu^2))^2 \right] \\ &\geq (\gamma + \kappa)^2 \cdot \mathbb{E}_{(s, a) \sim \mathcal{D}_E} \left[ \left( \operatorname{softmax}_{\underline{a} \in \mathcal{A}}^\beta Q(s, \underline{a}; \nu^1) - \operatorname{softmax}_{\underline{a} \in \mathcal{A}}^\beta Q(s, \underline{a}; \nu^2) \right)^2 \right]. \end{aligned}$$

The following proposition parallels Theorem B.2, which characterizes the global optimality and convergence of soft Q-learning. Recall that  $\nu_t$  is the PDE solution in (B.8) and  $\theta^{(m)}(k)$  is the soft Q-learning dynamics in (B.7).

**Proposition B.4.** There exists a unique fixed point solution to the projected soft Bellman optimality equation  $Q = \Pi_{\mathcal{F}} \mathcal{T}_\beta Q$ , which takes the form of  $Q^\dagger(x) = \int \sigma(x; \theta) d\underline{\nu}(\theta)$ . We assume that  $D_{\chi^2}(\underline{\nu} \parallel \nu_0) < \infty$  and  $\underline{\nu}(\theta) > 0$  for any  $\theta \in \mathbb{R}^D$ . Under Assumptions 4.1, 4.2, and B.3, it holds for  $\eta = \alpha^{-2}$  that

$$\inf_{t \in [0, T]} \mathbb{E}_{x \sim \mathcal{D}_E} \left[ (Q(x; \nu_t) - Q^\dagger(x))^2 \right] \leq \frac{(\kappa + \gamma) \cdot D_{\chi^2}(\underline{\nu} \parallel \nu_0)}{2\kappa \cdot T} + \frac{(\kappa + \gamma) \cdot C_*}{\kappa \cdot \alpha},$$

where  $C_* > 0$  is a constant depending on  $D_{\chi^2}(\underline{\nu} \parallel \nu_0)$ ,  $B_1$ ,  $B_2$ , and  $B_r$ . Moreover, it holds with probability at least  $1 - \delta$  that

$$\min_{\substack{k \leq T/\epsilon \\ (k \in \mathbb{N})}} \mathbb{E}_{x \sim \mathcal{D}_E} \left[ \left( \widehat{Q}(x; \theta^{(m)}(k)) - Q^\dagger(x) \right)^2 \right] \leq \frac{(\kappa + \gamma) \cdot D_{\chi^2}(\underline{\nu} \parallel \nu_0)}{2\kappa \cdot T} + \frac{(\kappa + \gamma) \cdot C_*}{\kappa \cdot \alpha} + \Delta(\epsilon, m, \delta, T),$$

where  $\Delta(\epsilon, m, \delta, T) > 0$  is an error term such that

$$\lim_{m \rightarrow \infty} \lim_{\epsilon \rightarrow 0^+} \Delta(\epsilon, m, \delta, T) = 0.$$

*Proof.* Replacing the max operator by the softmax operator in the proof of Theorem B.2 in §B.4 implies Proposition B.4.  $\square$

Moreover, soft Q-learning is equivalent to a variant of policy gradient [37, 57, 58, 61]. Hence, Proposition B.4 also characterizes the global optimality and convergence of such a variant of policy gradient.

#### B.4 Proof of Theorem B.2

For notational simplicity, we denote by  $\mathbb{E}_{\mathcal{D}_E}$  the expectation with respect to  $x \sim \mathcal{D}_E$  and  $\mathbb{E}_{\tilde{\mathcal{D}}_E}$  the expectation with respect to  $(x, r, x') \sim \tilde{\mathcal{D}}_E$ .

*Proof.* In parallel to the proof of Lemma 5.1 in §C.1, to establish the existence and uniqueness of the fixed point solution to the projected Bellman optimality equation  $Q = \Pi_{\mathcal{F}} \mathcal{T}^* Q$ , it suffices to show that  $\Pi_{\mathcal{F}} \mathcal{T}^* : \mathcal{F} \rightarrow \mathcal{F}$  is a contraction mapping. In particular, it holds for any  $Q^1, Q^2 \in \mathcal{F}$  that

$$\begin{aligned} \|\Pi_{\mathcal{F}} \mathcal{T}^* Q^1 - \Pi_{\mathcal{F}} \mathcal{T}^* Q^2\|_{\mathcal{L}_2(\mathcal{D}_E)}^2 &\leq \gamma^2 \cdot \mathbb{E}_{\tilde{\mathcal{D}}_E} \left[ \left( \max_{a \in \mathcal{A}} Q^1(s', \underline{a}) - \max_{a \in \mathcal{A}} Q^2(s', \underline{a}) \right)^2 \right] \\ &= \gamma^2 \cdot \mathbb{E}_{\mathcal{D}_E} \left[ \left( \max_{a \in \mathcal{A}} Q^1(s, \underline{a}) - \max_{a \in \mathcal{A}} Q^2(s, \underline{a}) \right)^2 \right] \\ &\leq \frac{\gamma^2}{(\gamma + \kappa)^2} \cdot \mathbb{E}_{\mathcal{D}_E} \left[ \left( Q^1(s, a) - Q^2(s, a) \right)^2 \right], \end{aligned}$$

where the equality follows from the fact that  $\mathcal{D}_E$  is the stationary distribution and the last inequality follows from Assumption B.1. Thus,  $\Pi_{\mathcal{F}} \mathcal{T}^* : \mathcal{F} \rightarrow \mathcal{F}$  is a contraction mapping. Following from the Banach fixed point theorem [28], there exists a unique fixed point solution  $Q^\dagger \in \mathcal{F}$  to the projected Bellman optimality equation  $Q = \Pi_{\mathcal{F}} \mathcal{T}^* Q$ . Moreover, in parallel to the proof of Lemma 5.1 in §C.1, there exists  $\nu^\dagger \in \mathcal{P}_2(\mathbb{R}^D)$  such that  $Q(x; \nu^\dagger) = Q^\dagger(x)$ ,  $h(x; \nu^\dagger) = 0$ , and  $\mathcal{W}_2(\nu^\dagger, \nu_0) \leq \alpha^{-1} \cdot \bar{D}$ , where  $\bar{D} = D_{\chi^2}(\bar{\nu} \parallel \nu_0)^{1/2}$ .

For notational simplicity, we define  $Q^{\mathcal{A}}(x) = \max_{\underline{a} \in \mathcal{A}} Q(s, \underline{a})$ . In parallel to (C.13) in the proof of Lemma 5.2 in §C.2, we have that

$$\frac{d}{dt} \frac{\mathcal{W}_2(\nu_t, \nu^\dagger)^2}{2} = \underbrace{\eta \cdot \int_0^1 \langle \partial_s h(\cdot; \beta_s), v_s \rangle_{\beta_s} ds}_{(i)} + \underbrace{\eta \cdot \int_0^1 \int \langle h(\theta; \beta_s), \partial_s(v_s \cdot \beta_s)(\theta) \rangle d\theta ds}_{(ii)}, \quad (\text{B.9})$$

where  $\beta : [0, 1] \rightarrow \mathcal{P}_2(\mathbb{R}^D)$  is the geodesic connecting  $\nu_t$  and  $\nu^\dagger$  with  $\partial_s \beta_s = -\text{div}(\beta_s \cdot v_s)$ .

**Upper bounding term (i) of (B.9).** In parallel to (C.5) and (C.6) in the proof of Lemma C.1, we have that

$$\begin{aligned} \langle \partial_s h(\cdot; \beta_s), v_s \rangle_{\beta_s} &= -\mathbb{E}_{\tilde{\mathcal{D}}_E} \left[ \partial_s (Q(x; \beta_s) - \gamma \cdot Q^{\mathcal{A}}(x'; \beta_s)) \cdot \partial_s Q(x; \beta_s) \right] \\ &\leq -\mathbb{E}_{\mathcal{D}_E} \left[ (\partial_s Q(x; \beta_s))^2 \right] + \gamma \cdot \mathbb{E}_{\mathcal{D}_E} \left[ (\partial_s Q(x; \beta_s))^2 \right]^{1/2} \cdot \mathbb{E}_{\mathcal{D}_E} \left[ (\partial_s Q^{\mathcal{A}}(x; \beta_s))^2 \right]^{1/2}. \end{aligned} \quad (\text{B.10})$$

For the second term on the right-hand side of (B.10), we have that

$$\begin{aligned} \mathbb{E}_{\mathcal{D}_E} \left[ (\partial_s Q^{\mathcal{A}}(x; \beta_s))^2 \right] &= \lim_{u \rightarrow 0} \mathbb{E}_{\mathcal{D}_E} \left[ \left( u^{-1} \cdot (Q^{\mathcal{A}}(x; \beta_{s+u}) - Q^{\mathcal{A}}(x; \beta_s)) \right)^2 \right] \\ &\leq (\gamma + \kappa)^{-2} \cdot \lim_{u \rightarrow 0} u^{-2} \cdot \mathbb{E}_{\mathcal{D}_E} \left[ (Q(x; \beta_{s+u}) - Q(x; \beta_s))^2 \right] \\ &= (\gamma + \kappa)^{-2} \cdot \mathbb{E}_{\mathcal{D}_E} \left[ (\partial_s Q(x; \beta_s))^2 \right], \end{aligned} \quad (\text{B.11})$$

where the inequality follows from Assumption B.1 and the fact that  $Q(\cdot; \nu) \in \alpha \cdot \mathcal{F}$ . Plugging (B.11) into (B.10), we have that

$$\langle \partial_s h(\cdot; \beta_s), v_s \rangle_{\beta_s} \leq -\frac{\kappa}{\gamma + \kappa} \cdot \mathbb{E}_{\mathcal{D}_E} \left[ (\partial_s Q(x; \beta_s))^2 \right],$$

which further implies that

$$\begin{aligned} \int_0^1 \langle \partial_s h(\cdot; \beta_s), v_s \rangle_{\beta_s} ds &\leq -\frac{\kappa}{\gamma + \kappa} \cdot \int_0^1 \mathbb{E}_{\mathcal{D}_E} \left[ (\partial_s Q(x; \beta_s))^2 \right] ds \\ &\leq -\frac{\kappa}{\gamma + \kappa} \cdot \mathbb{E}_{\mathcal{D}_E} \left[ \left( \int_0^1 \partial_s Q(x; \beta_s) ds \right)^2 \right] \\ &= -\frac{\kappa}{\gamma + \kappa} \cdot \mathbb{E}_{\mathcal{D}_E} \left[ (Q(x; \nu_t) - Q(x; \nu^\dagger))^2 \right]. \end{aligned} \quad (\text{B.12})$$

**Upper bounding term (ii) of (B.9).** In parallel to the proof of Lemma C.2 in §C.2, noting that  $|Q^A(x; \nu)| \leq \sup_{x \in \mathcal{X}} |Q(x; \nu)|$  for any  $\nu \in \mathcal{P}_2(\mathbb{R}^D)$ , we have that

$$\|\nabla_\theta h(\theta; \nu_t)\|_{\mathbb{F}} \leq \alpha \cdot B_2 \cdot (2\alpha \cdot B_1 \cdot \mathcal{W}_2(\nu_t, \nu_0) + B_r).$$

In parallel to (C.15) and (C.16), we have that

$$\int_0^1 \int \left| \langle h(\theta; \beta_s), \partial_s(v_s \cdot \beta_s)(\theta) \rangle \right| d\theta ds \leq C_* \cdot \alpha^{-1}, \quad (\text{B.13})$$

where  $C_* > 0$  is a constant that depends on  $\bar{D}$ ,  $B_1$ ,  $B_2$ , and  $B_r$ .

Plugging (B.12) and (B.13) into (B.9), we have that

$$\frac{d}{dt} \frac{\mathcal{W}_2(\nu_t, \nu^\dagger)^2}{2} \leq -\frac{\eta \cdot \kappa}{\gamma + \kappa} \cdot \mathbb{E}_{\mathcal{D}_E} \left[ (Q(x; \nu_t) - Q(x; \nu^\dagger))^2 \right] + C_* \cdot \eta \cdot \alpha^{-1}.$$

Thus, in parallel to the proof of Theorem 4.3 in §5, we have that

$$\inf_{t \in [0, T]} \mathbb{E}_{\mathcal{D}} \left[ (Q(x; \nu_t) - Q^\dagger(x))^2 \right] \leq \frac{(\kappa + \gamma) \cdot D_{\chi^2}(\bar{\nu} \parallel \nu_0)}{2\kappa \cdot T} + C_* \cdot \alpha^{-1} \cdot \frac{\kappa + \gamma}{\kappa},$$

which completes the proof of (B.5) in Theorem B.2. Meanwhile, in parallel to the proof of Lemma D.6 in §D.2, we upper bound the error of approximating  $\hat{\nu}_k$  by  $\nu_{k\epsilon}$ , which further implies (B.6) of Theorem B.2.  $\square$

## C Proofs of Supporting Lemmas

For notational simplicity, we denote by  $\mathbb{E}_{\mathcal{D}}$  the expectation with respect to  $x \sim \mathcal{D}$  and  $\mathbb{E}_{\tilde{\mathcal{D}}}$  the expectation with respect to  $(x, r, x') \sim \tilde{\mathcal{D}}$ . Also, with a slight abuse of notations, we write  $\theta^{(m)} = \{\theta_i\}_{i=1}^m$ .

### C.1 Proof of Lemma 5.1

*Proof. Existence and uniqueness of  $Q^*$ .* To establish the existence of the fixed point solution  $Q^*$  to the projected Bellman equation  $Q = \Pi_{\mathcal{F}} \mathcal{T}^\pi Q$ , it suffices to show that  $\Pi_{\mathcal{F}} \mathcal{T}^\pi : \mathcal{F} \rightarrow \mathcal{F}$  is a contraction mapping. It holds for any  $Q^1, Q^2 \in \mathcal{F}$  that

$$\begin{aligned} \|\Pi_{\mathcal{F}} \mathcal{T}^\pi Q^1 - \Pi_{\mathcal{F}} \mathcal{T}^\pi Q^2\|_{\mathcal{L}_2(\mathcal{D})}^2 &\leq \gamma^2 \cdot \mathbb{E}_{\tilde{\mathcal{D}}} \left[ (Q^1(x') - Q^2(x'))^2 \right] \\ &= \gamma^2 \cdot \|Q^1 - Q^2\|_{\mathcal{L}_2(\mathcal{D})}^2, \end{aligned}$$

where the last equality follows from the fact that  $\mathcal{D}$  is the stationary distribution. Thus,  $\Pi_{\mathcal{F}} \mathcal{T}^\pi : \mathcal{F} \rightarrow \mathcal{F}$  is a contraction mapping. Note that  $\mathcal{F}$  is complete. Following from the Banach fixed point theorem [28], there exists a unique  $Q^* \in \mathcal{F}$  that solves the projected Bellman equation  $Q = \Pi_{\mathcal{F}} \mathcal{T}^\pi Q$ . Moreover, by the definition of  $\mathcal{F}$  in (4.3), there exists  $\bar{\rho} \in \mathcal{P}_2(\mathbb{R}^D)$  such that

$$Q^*(x) = \int \sigma(x; \theta) d\bar{\rho}(\theta).$$

**Proof of (i) in Lemma 5.1.** We define

$$\rho^* = \rho_0 + \alpha^{-1} \cdot (\bar{\rho} - \rho_0). \quad (\text{C.1})$$

By the definition of  $Q(\cdot; \rho)$  in (3.2) and the fact that  $Q(x; \rho_0) = 0$ , we have that  $Q(x; \rho^*) = Q^*(x)$ , which completes the proof of (i) in Lemma 5.1.

**Proof of (ii) in Lemma 5.1.** For (ii) of Lemma 5.1, note that  $Q(\cdot; \rho^*) = \Pi_{\mathcal{F}} \mathcal{T}^\pi Q(\cdot; \rho^*)$ . Thus, we have that

$$\langle Q(\cdot; \rho^*) - \mathcal{T}^\pi Q(\cdot; \rho^*), f(\cdot) - Q(\cdot; \rho^*) \rangle_{\mathcal{D}} \geq 0, \quad \forall f \in \mathcal{F},$$

which further implies that

$$\mathbb{E}_{\bar{\mathcal{D}}}\left[(Q(x; \rho^*) - r - \gamma \cdot Q(x'; \rho^*)) \cdot \int \sigma(x; \theta) d(\rho - \bar{\rho})(\theta)\right] \geq 0, \quad \forall \rho \in \mathcal{P}_2(\mathbb{R}^D). \quad (\text{C.2})$$

Let  $\rho = (\text{id} + h \cdot v)_\# \bar{\rho}$  for a sufficiently small scaling parameter  $h \in \mathbb{R}_+$  and any Lipschitz-continuous mapping  $v : \mathbb{R}^D \rightarrow \mathbb{R}^D$ . Then, following from (C.2), we have that

$$\int \mathbb{E}_{\bar{\mathcal{D}}}\left[(Q(x; \rho^*) - r - \gamma \cdot Q(x'; \rho^*)) \cdot (\sigma(x; \theta + h \cdot v(\theta)) - \sigma(x; \theta))\right] d\bar{\rho}(\theta) \geq 0 \quad (\text{C.3})$$

for any  $v : \mathbb{R}^D \rightarrow \mathbb{R}^D$ . Dividing the both sides of (C.3) by  $h$  and letting  $h \rightarrow 0^+$ , we have for any  $v : \mathbb{R}^D \rightarrow \mathbb{R}^D$  that

$$\begin{aligned} 0 &\leq \int \mathbb{E}_{\bar{\mathcal{D}}}\left[(Q(x; \rho^*) - r - \gamma \cdot Q(x'; \rho^*)) \cdot \langle \nabla_\theta \sigma(x; \theta), v(\theta) \rangle\right] d\bar{\rho}(\theta) \\ &= -\alpha^{-1} \cdot \int \langle g(\theta; \rho^*), v(\theta) \rangle d\bar{\rho}(\theta), \end{aligned}$$

where the equality follows from the definition of  $g$  in (3.5). Thus, we have that  $g(\theta; \rho^*) = 0$  for  $\bar{\rho}$ -a.e., which completes the proof of (ii) in Lemma 5.1.

**Proof of (iii) in Lemma 5.1.** Following from the definition of  $\rho^*$  in (C.1), we have that

$$\begin{aligned} D_{\chi^2}(\rho^* \parallel \rho_0) &= \int \left( \frac{\rho^*(\theta)}{\rho_0(\theta)} - 1 \right)^2 d\rho_0(\theta) = \int \left( \frac{(1 - \alpha^{-1}) \cdot \rho_0(\theta) + \alpha^{-1} \cdot \bar{\rho}(\theta)}{\rho_0(\theta)} - 1 \right)^2 d\rho_0(\theta) = \alpha^{-2} \cdot \bar{D}^2, \end{aligned}$$

where  $\bar{D} = D_{\chi^2}(\bar{\rho} \parallel \rho_0)^{1/2}$ . By Lemma E.3, we have that

$$\mathcal{W}_2(\rho^*, \rho_0) \leq D_{\text{KL}}(\rho^* \parallel \rho_0)^{1/2} \leq D_{\chi^2}(\rho^* \parallel \rho_0)^{1/2} \leq \alpha^{-1} \cdot \bar{D},$$

which completes the proof of (iii) in Lemma 5.1.  $\square$

## C.2 Proof of Lemma 5.2

We first introduce the following lemmas. The first lemma establishes the one-point monotonicity of  $g(\cdot; \beta_t)$  along a curve  $\beta : [0, 1] \rightarrow \mathcal{P}_2(\mathbb{R}^D)$  on the Wasserstein space.

**Lemma C.1.** Let  $\beta : [0, 1] \rightarrow \mathcal{P}_2(\mathbb{R}^D)$  be a curve such that  $\partial_t \beta_t = -\text{div}(\beta_t \cdot v_t)$  for a vector field  $v$ . We have that

$$\langle \partial_t g(\cdot; \beta_t), v_t \rangle_{\beta_t} \leq -(1 - \gamma) \cdot \mathbb{E}_{\mathcal{D}}\left[(\partial_t Q(x; \beta_t))^2\right].$$

Furthermore, we have that

$$\int_0^1 \langle \partial_s g(\cdot; \beta_s), v_s \rangle_{\beta_s} ds \leq -(1 - \gamma) \cdot \mathbb{E}_{\mathcal{D}}\left[(Q(x; \beta_0) - Q(x; \beta_1))^2\right]. \quad (\text{C.4})$$

*Proof.* Following from the definition of  $g$  in (3.5), we have that

$$\partial_t g(\theta; \beta_t) = -\alpha \cdot \mathbb{E}_{\bar{\mathcal{D}}}\left[\partial_t (Q(x; \beta_t) - \gamma \cdot Q(x'; \beta_t)) \cdot \nabla_\theta \sigma(x; \theta)\right].$$

Thus, following from integration by parts and the continuity equation  $\partial_t \beta_t = -\text{div}(\beta_t \cdot v_t)$ , we have that

$$\begin{aligned} \langle \partial_t g(\cdot; \beta_t), v_t \rangle_{\beta_t} &= -\int \left\langle \alpha \cdot \mathbb{E}_{\bar{\mathcal{D}}}\left[\partial_t (Q(x; \beta_t) - \gamma \cdot Q(x'; \beta_t)) \cdot \nabla_\theta \sigma(x; \theta)\right], v_t(\theta) \cdot \beta_t(\theta) \right\rangle d\theta \\ &= -\int \alpha \cdot \mathbb{E}_{\bar{\mathcal{D}}}\left[\partial_t (Q(x; \beta_t) - \gamma \cdot Q(x'; \beta_t)) \cdot \sigma(x; \theta)\right] \cdot \partial_t \beta_t(\theta) d\theta \\ &= -\mathbb{E}_{\bar{\mathcal{D}}}\left[\partial_t (Q(x; \beta_t) - \gamma \cdot Q(x'; \beta_t)) \cdot \partial_t Q(x; \beta_t)\right], \end{aligned} \quad (\text{C.5})$$

where the last equality follows from the definition of  $Q$  in (3.2). Applying the Cauchy-Schwartz inequality to (C.5), we have that

$$\begin{aligned} \langle \partial_t g(\cdot; \beta_t), v_t \rangle_{\beta_t} &= -\mathbb{E}_{\tilde{\mathcal{D}}} \left[ (\partial_t Q(x; \beta_t))^2 \right] + \gamma \cdot \mathbb{E}_{\tilde{\mathcal{D}}} [\partial_t Q(x'; \beta_t) \cdot \partial_t Q(x; \beta_t)] \\ &\leq -\mathbb{E}_{\tilde{\mathcal{D}}} \left[ (\partial_t Q(x; \beta_t))^2 \right] + \gamma \cdot \mathbb{E}_{\tilde{\mathcal{D}}} \left[ (\partial_t Q(x; \beta_t))^2 \right]^{1/2} \cdot \mathbb{E}_{\tilde{\mathcal{D}}} \left[ (\partial_t Q(x'; \beta_t))^2 \right]^{1/2} \\ &= -(1 - \gamma) \cdot \mathbb{E}_{\mathcal{D}} \left[ (\partial_t Q(x; \beta_t))^2 \right], \end{aligned} \quad (\text{C.6})$$

where the last equality follows from the fact that the marginal distributions of  $\tilde{\mathcal{D}}$  with respect to  $x$  and  $x'$  are  $\mathcal{D}$ , since  $\mathcal{D}$  is the stationary distribution. Furthermore, we have that

$$\begin{aligned} \int_0^1 \langle \partial_s g(\cdot; \beta_s), v_s \rangle_{\beta_s} ds &\leq -(1 - \gamma) \cdot \int_0^1 \mathbb{E}_{\mathcal{D}} \left[ (\partial_s Q(x; \beta_s))^2 \right] ds \\ &\leq -(1 - \gamma) \cdot \mathbb{E}_{\mathcal{D}} \left[ \left( \int_0^1 \partial_s Q(x; \beta_s) ds \right)^2 \right] \\ &= -(1 - \gamma) \cdot \mathbb{E}_{\mathcal{D}} \left[ (Q(x; \beta_1) - Q(x; \beta_0))^2 \right], \end{aligned}$$

which completes the proof of Lemma C.1.  $\square$

The following lemma upper bounds the norms of  $Q$  and  $\nabla_\theta g$ .

**Lemma C.2.** Under Assumptions 4.1 and 4.2, it holds for any  $\rho \in \mathcal{P}_2(\mathbb{R}^D)$  that

$$\sup_{x \in \mathcal{X}} |Q(x; \rho)| \leq \alpha \cdot \min\{B_1 \cdot \mathcal{W}_2(\rho, \rho_0), B_0\}, \quad (\text{C.7})$$

$$\sup_{\theta \in \mathbb{R}^D} \|\nabla_\theta g(\theta; \rho)\|_{\text{F}} \leq \alpha \cdot B_2 \cdot \min\{2\alpha \cdot B_1 \cdot \mathcal{W}_2(\rho, \rho_0) + B_r, 2\alpha \cdot B_0 + B_r\}. \quad (\text{C.8})$$

*Proof.* We introduce the Wasserstein-1 distance, which is defined as

$$\mathcal{W}_1(\mu^1, \mu^2) = \inf \left\{ \mathbb{E}[\|X - Y\|] \mid \text{law}(X) = \mu^1, \text{law}(Y) = \mu^2 \right\}$$

for any  $\mu^1, \mu^2 \in \mathcal{P}(\mathbb{R}^D)$  with finite first moments. Thus, we have that  $\mathcal{W}_1(\mu^1, \mu^2) \leq \mathcal{W}_2(\mu^1, \mu^2)$ . The Wasserstein-1 distance has the following dual representation [5],

$$\mathcal{W}_1(\mu^1, \mu^2) = \sup \left\{ \int f(x) d(\mu^1 - \mu^2)(x) \mid \text{continuous } f : \mathbb{R}^D \rightarrow \mathbb{R}, \text{Lip}(f) \leq 1 \right\}. \quad (\text{C.9})$$

Following from Assumptions 4.1 and 4.2, we have that  $\|\nabla_\theta \sigma(x; \theta)\| \leq B_1$  for any  $x \in \mathcal{X}$  and  $\theta \in \mathbb{R}^D$ , which implies that  $\text{Lip}(\sigma(x; \cdot)/B_1) \leq 1$  for any  $x \in \mathcal{X}$ . Note that  $Q(x; \rho_0) = 0$  for any  $x \in \mathcal{X}$ . Thus, by (C.9) we have for any  $\rho \in \mathcal{P}_2(\mathbb{R}^D)$  and  $x \in \mathcal{X}$  that

$$|Q(x; \rho)| = \alpha \cdot \left| \int \sigma(x; \theta) \cdot d(\rho - \rho_0)(\theta) \right| \leq \alpha \cdot B_1 \cdot \mathcal{W}_1(\rho, \rho_0) \leq \alpha \cdot B_1 \cdot \mathcal{W}_2(\rho, \rho_0). \quad (\text{C.10})$$

Meanwhile, following from Assumptions 4.1 and 4.2, we have for any  $x \in \mathcal{X}$  and  $\rho \in \mathcal{P}_2(\mathbb{R}^D)$  that

$$|Q(x; \rho)| = \alpha \cdot \left| \int \sigma(x; \theta) d\rho(\theta) \right| \leq \alpha \cdot B_0. \quad (\text{C.11})$$

Combining (C.10) and (C.11), we have for any  $\rho \in \mathcal{P}_2(\mathbb{R}^D)$  that

$$\sup_{x \in \mathcal{X}} |Q(x; \rho)| \leq \alpha \cdot \min\{B_1 \cdot \mathcal{W}_2(\rho, \rho_0), B_0\}, \quad (\text{C.12})$$

which completes the proof of (C.7) in Lemma C.2. Following from the definition of  $g$  in (3.5), we have for any  $x \in \mathcal{X}$  and  $\rho \in \mathcal{P}_2(\mathbb{R}^D)$  that

$$\begin{aligned} \|\nabla_\theta g(\theta; \rho)\|_{\text{F}} &\leq \alpha \cdot \mathbb{E}_{\tilde{\mathcal{D}}} \left[ |Q(x; \rho) - r - \gamma \cdot Q(x'; \rho)| \cdot \|\nabla_{\theta\theta}^2 \sigma(x; \theta)\|_{\text{F}} \right] \\ &\leq \alpha \cdot \min\{2\alpha \cdot B_1 \cdot \mathcal{W}_2(\rho, \rho_0) + B_r, 2\alpha \cdot B_0 + B_r\} \cdot B_2. \end{aligned}$$

Here the last inequality follows from (C.12) and the fact that  $\|\nabla_{\theta\theta}^2 \sigma(x; \theta)\|_{\text{F}} \leq B_2$  for any  $x \in \mathcal{X}$  and  $\rho \in \mathcal{P}_2(\mathbb{R}^D)$ , which follows from Assumptions 4.1 and 4.2. Thus, we complete the proof of Lemma C.2.  $\square$

We are now ready to present the proof of Lemma 5.2.

*Proof.* Recall that  $\rho_t$  is the PDE solution in (3.4), that is,

$$\partial_t \rho_t = -\eta \cdot \operatorname{div}(\rho_t \cdot g(\cdot; \rho_t)),$$

where

$$g(\theta; \rho) = -\alpha \cdot \mathbb{E}_{\bar{\mathcal{D}}} \left[ (Q(x; \rho) - r - \gamma \cdot Q(x'; \rho)) \cdot \nabla_{\theta} \sigma(x; \theta) \right].$$

We fix a  $t \in [0, T]$ . We denote by  $\beta : [0, 1] \rightarrow \mathcal{P}_2(\mathbb{R}^D)$  the geodesic connecting  $\rho_t$  and  $\rho^*$ . Specifically,  $\beta$  satisfies that  $\beta'_s = -\operatorname{div}(\beta_s \cdot v_s)$  for a vector field  $v$ . Following from Lemma E.2, we have that

$$\begin{aligned} \frac{d}{dt} \frac{\mathcal{W}_2(\rho_t, \rho^*)^2}{2} &= -\eta \cdot \langle g(\cdot; \rho_t), v_0 \rangle_{\rho_t} \\ &= \eta \cdot \int_0^1 \partial_s \langle g(\cdot; \beta_s), v_s \rangle_{\beta_s} ds - \eta \cdot \langle g(\cdot; \rho^*), v_1 \rangle_{\rho^*} \\ &= \eta \cdot \underbrace{\int_0^1 \langle \partial_s g(\cdot; \beta_s), v_s \rangle_{\beta_s} ds}_{(i)} + \eta \cdot \underbrace{\int_0^1 \int \langle g(\theta; \beta_s), \partial_s (v_s \cdot \beta_s)(\theta) \rangle d\theta ds}_{(ii)}, \end{aligned} \quad (\text{C.13})$$

where the last equality follows from (ii) of Lemma 5.1.

For term (i) of (C.13), following from (C.4) of Lemma C.1, we have that

$$\begin{aligned} \int_0^1 \langle \partial_s g(\cdot; \beta_s), v_s \rangle_{\beta_s} ds &\leq -(1 - \gamma) \cdot \mathbb{E}_{\mathcal{D}} \left[ (Q(x; \beta_0) - Q(x; \beta_1))^2 \right] \\ &= -(1 - \gamma) \cdot \mathbb{E}_{\mathcal{D}} \left[ (Q(x; \rho_t) - Q^*(x))^2 \right]. \end{aligned} \quad (\text{C.14})$$

For term (ii) of (C.14), we have that

$$\begin{aligned} \int \left| \langle g(\theta; \beta_s), \partial_s (v_s \cdot \beta_s)(\theta) \rangle \right| d\theta &= \int \left| \langle \nabla_{\theta} g(\theta; \beta_s), \beta_s(\theta) \cdot v_s(\theta) \otimes v_s(\theta) \rangle \right| d\theta \\ &\leq \sup_{\theta \in \mathbb{R}^D} \left\| \nabla_{\theta} g(\theta; \beta_s) \right\|_{\mathbb{F}} \cdot \|v_s\|_{\beta_s}^2, \end{aligned}$$

where the equality follows from integration by parts and Lemma E.4. Since  $\beta$  is the geodesic connecting  $\rho_t$  and  $\rho^*$ , (2.7) implies that  $\|v_s\|_{\beta_s}^2 = \mathcal{W}_2(\beta_0, \beta_1)^2 = \mathcal{W}_2(\rho_t, \rho^*)^2$  for any  $s \in [0, 1]$ . Applying (C.8) of Lemma C.2, we have that

$$\begin{aligned} \int \left| \langle g(\theta; \beta_s), \partial_s (v_s \cdot \beta_s)(\theta) \rangle \right| d\theta &\leq \alpha \cdot B_2 \cdot (2\alpha \cdot B_1 \cdot \mathcal{W}_2(\rho_t, \rho_0) + B_r) \cdot \mathcal{W}_2(\rho_t, \rho^*)^2 \\ &\leq 4\alpha \cdot B_2 \cdot (6\alpha \cdot B_1 \cdot \mathcal{W}_2(\rho_0, \rho^*) + B_r) \cdot \mathcal{W}_2(\rho_0, \rho^*)^2, \end{aligned} \quad (\text{C.15})$$

where the last inequality follows from the condition of Lemma 5.2 that  $\mathcal{W}_2(\rho_t, \rho^*) \leq 2\mathcal{W}_2(\rho_0, \rho^*)$  and the fact that  $\mathcal{W}_2(\rho_t, \rho_0) \leq \mathcal{W}_2(\rho_t, \rho^*) + \mathcal{W}_2(\rho_0, \rho^*)$ . Then, applying (iii) of Lemma 5.1 to (C.15), we have that

$$\begin{aligned} \int_0^1 \int \left| \langle g(\theta; \beta_s), \partial_s (v_s \cdot \beta_s)(\theta) \rangle \right| d\theta ds &\leq 4\alpha^{-1} \cdot B_2 \cdot \bar{D}^2 \cdot (6B_1 \cdot \bar{D} + B_r) \\ &= C_* \cdot \alpha^{-1}, \end{aligned} \quad (\text{C.16})$$

where  $C_* > 0$  is a constant depending on  $\bar{D}$ ,  $B_1$ ,  $B_2$ , and  $B_r$ .

Finally, plugging (C.14) and (C.16) into (C.13), we have that

$$\frac{d}{dt} \frac{\mathcal{W}_2(\rho_t, \rho^*)^2}{2} \leq -(1 - \gamma) \cdot \eta \cdot \mathbb{E}_{\mathcal{D}} \left[ (Q(x; \rho_t) - Q^*(x))^2 \right] + C_* \cdot \alpha^{-1} \cdot \eta,$$

which completes the proof of Lemma 5.2.  $\square$



## D Mean-Field Limit of Neural Networks

In this section, we prove Proposition 3.1, whose formal version is presented as follows. Recall that  $\rho_t$  is the PDE solution in (3.4) and  $\hat{\rho}_k = m^{-1} \cdot \sum_{i=1}^m \theta_i(k)$  is the empirical distribution of  $\theta^{(m)}(k) = \{\theta_i(k)\}_{i=1}^m$ . Note that we omit the dependence of  $\hat{\rho}_k$  on  $m$  and  $\epsilon$  for notational simplicity.

**Proposition D.1** (Formal Version of Proposition 3.1). Let  $f : \mathbb{R}^D \rightarrow \mathbb{R}$  be any continuous function such that  $\|f\|_\infty \leq 1$  and  $\text{Lip}(f) \leq 1$ . Under Assumptions 4.1 and 4.2, it holds that

$$\begin{aligned} & \sup_{\substack{k \leq T/\epsilon \\ (k \in \mathbb{N})}} \left| \int f(\theta) d\rho_{k\epsilon}(\theta) - \int f(\theta) d\hat{\rho}_k(\theta) \right| \\ & \leq B \cdot e^{BT} \cdot \left( \sqrt{\log(m/\delta)/m} + \sqrt{\epsilon \cdot (D + \log(m/\delta))} \right) \end{aligned}$$

with probability at least  $1 - \delta$ . Here  $B$  is a constant that depends on  $\alpha, \eta, \gamma, B_r$ , and  $B_j$  ( $j \in \{0, 1, 2\}$ ).

The proof of Proposition D.1 is based on [6, 53, 54], which utilizes the propagation of chaos [66]. Recall that  $g(\cdot; \rho)$  is a vector field defined as follows,

$$g(\theta; \rho) = -\alpha \cdot \mathbb{E}_{\tilde{\mathcal{D}}} \left[ (Q(x; \rho) - r - \gamma \cdot Q(x'; \rho)) \cdot \nabla_\theta \sigma(x; \theta) \right].$$

Correspondingly, we define the finite-width and stochastic counterparts of  $g(\theta; \rho)$  as follows,

$$\hat{g}(\theta; \theta^{(m)}) = -\alpha \cdot \mathbb{E}_{\tilde{\mathcal{D}}} \left[ (\hat{Q}(x; \theta^{(m)}) - r - \gamma \cdot \hat{Q}(x'; \theta^{(m)})) \cdot \nabla_\theta \sigma(x; \theta) \right], \quad (\text{D.1})$$

$$\hat{G}_k(\theta; \theta^{(m)}) = -\alpha \cdot (\hat{Q}(x_k; \theta^{(m)}) - r_k - \gamma \cdot \hat{Q}(x'_k; \theta^{(m)})) \cdot \nabla_\theta \sigma(x_k; \theta), \quad (\text{D.2})$$

where  $(x_k, r_k, x'_k) \sim \tilde{\mathcal{D}}$ . Following from [6, 53], we consider the following four dynamics.

- **Temporal-difference (TD).** We consider the following TD dynamics  $\theta^{(m)}(k)$ , where  $k \in \mathbb{N}$ , with  $\theta_i(0) \stackrel{\text{i.i.d.}}{\sim} \rho_0$  ( $i \in [m]$ ) as its initialization,

$$\begin{aligned} \theta_i(k+1) &= \theta_i(k) - \eta\epsilon \cdot \alpha \cdot \left( \hat{Q}(x_k; \theta^{(m)}(k)) - r_k - \gamma \cdot \hat{Q}(x'_k; \theta^{(m)}(k)) \right) \cdot \nabla_\theta \sigma(x_k; \theta_i(k)) \\ &= \theta_i(k) + \eta\epsilon \cdot \hat{G}_k(\theta_i(k); \theta^{(m)}(k)), \end{aligned} \quad (\text{D.3})$$

where  $(x_k, r_k, x'_k) \sim \tilde{\mathcal{D}}$ . Note that this definition is equivalent to (2.3).

- **Expected temporal-difference (ETD).** We consider the following expected TD dynamics  $\check{\theta}^{(m)}(k)$ , where  $k \in \mathbb{N}$ , with  $\check{\theta}_i(0) = \theta_i(0)$  ( $i \in [m]$ ) as its initialization,

$$\begin{aligned} \check{\theta}_i(k+1) &= \check{\theta}_i(k) - \eta\epsilon \cdot \alpha \cdot \mathbb{E}_{\tilde{\mathcal{D}}} \left[ \left( \hat{Q}(x; \check{\theta}^{(m)}(k)) - r - \gamma \cdot \hat{Q}(x'; \check{\theta}^{(m)}(k)) \right) \cdot \nabla_\theta \sigma(x; \check{\theta}_i(k)) \right] \\ &= \check{\theta}_i(k) + \eta\epsilon \cdot \hat{g}(\check{\theta}_i(k); \check{\theta}^{(m)}(k)). \end{aligned} \quad (\text{D.4})$$

- **Continuous-time temporal-difference (CTTD).** We consider the following continuous-time TD dynamics  $\tilde{\theta}^{(m)}(t)$ , where  $t \in \mathbb{R}_+$ , with  $\tilde{\theta}_i(0) = \theta_i(0)$  ( $i \in [m]$ ) as its initialization,

$$\begin{aligned} \frac{d}{dt} \tilde{\theta}_i(t) &= -\eta \cdot \alpha \cdot \mathbb{E}_{\tilde{\mathcal{D}}} \left[ \left( \hat{Q}(x; \tilde{\theta}^{(m)}(t)) - r - \gamma \cdot \hat{Q}(x'; \tilde{\theta}^{(m)}(t)) \right) \cdot \nabla_\theta \sigma(x; \tilde{\theta}_i(t)) \right] \\ &= \eta \cdot \hat{g}(\tilde{\theta}_i(t); \tilde{\theta}^{(m)}(t)). \end{aligned} \quad (\text{D.5})$$

- **Ideal particle (IP).** We consider the following ideal particle dynamics  $\bar{\theta}^{(m)}(t)$ , where  $t \in \mathbb{R}_+$ , with  $\bar{\theta}_i(0) = \theta_i(0)$  ( $i \in [m]$ ) as its initialization,

$$\begin{aligned} \frac{d}{dt} \bar{\theta}_i(t) &= -\eta \cdot \alpha \cdot \mathbb{E}_{\tilde{\mathcal{D}}} \left[ (Q(x; \rho_t) - r - \gamma \cdot Q(x'; \rho_t)) \cdot \nabla_\theta \sigma(x; \bar{\theta}_i(t)) \right] \\ &= \eta \cdot g(\bar{\theta}_i(t); \rho_t), \end{aligned} \quad (\text{D.6})$$

where  $\rho_t$  is the PDE solution in (3.4).

We aim to prove that  $\hat{\rho}_k = m^{-1} \cdot \sum_{i=1}^m \delta_{\theta_i(k)}$  weakly converges to  $\rho_{k\epsilon}$ . For any continuous function  $f : \mathbb{R}^D \rightarrow \mathbb{R}$  such that  $\|f\|_\infty \leq 1$  and  $\text{Lip}(f) \leq 1$ , we use the IP, CTTD, and ETD dynamics as the interpolating dynamics,

$$\begin{aligned}
& \overbrace{\left| \int f(\theta) d\rho_{k\epsilon}(\theta) - \int f(\theta) d\hat{\rho}_k(\theta) \right|}^{\text{PDE} - \text{TD}} \\
& \leq \left| \int f(\theta) d\rho_{k\epsilon}(\theta) - m^{-1} \cdot \sum_{i=1}^m f(\bar{\theta}_i(k\epsilon)) \right| + \left| m^{-1} \cdot \sum_{i=1}^m f(\bar{\theta}_i(k\epsilon)) - m^{-1} \cdot \sum_{i=1}^m f(\tilde{\theta}_i(k\epsilon)) \right| \\
& \quad + \left| m^{-1} \cdot \sum_{i=1}^m f(\tilde{\theta}_i(k\epsilon)) - m^{-1} \cdot \sum_{i=1}^m f(\check{\theta}_i(k)) \right| + \left| m^{-1} \cdot \sum_{i=1}^m f(\check{\theta}_i(k)) - m^{-1} \cdot \sum_{i=1}^m f(\theta_i(k)) \right| \\
& \leq \underbrace{\left| \int f(\theta) d\rho_{k\epsilon}(\theta) - m^{-1} \cdot \sum_{i=1}^m f(\bar{\theta}_i(k\epsilon)) \right|}_{\text{PDE} - \text{IP}} + \underbrace{\|\bar{\theta}^{(m)}(k\epsilon) - \tilde{\theta}^{(m)}(k\epsilon)\|_{(m)}}_{\text{IP} - \text{CTTD}} \\
& \quad + \underbrace{\|\tilde{\theta}^{(m)}(k\epsilon) - \check{\theta}^{(m)}(k)\|_{(m)}}_{\text{CTTD} - \text{ETD}} + \underbrace{\|\check{\theta}^{(m)}(k) - \theta^{(m)}(k)\|_{(m)}}_{\text{ETD} - \text{TD}}, \tag{D.7}
\end{aligned}$$

where the last inequality follows from the fact that  $\text{Lip}(f) \leq 1$ . Here the norm  $\|\cdot\|_{(m)}$  of  $\theta^{(m)} = \{\theta_i\}_{i=1}^m$  is defined as follows,

$$\|\theta^{(m)}\|_{(m)} = \sup_{i \in [m]} \|\theta_i\|. \tag{D.8}$$

In what follows, we define  $B > 0$  as a constant that depends on  $\alpha, \eta, \gamma, B_r$ , and  $B_j$  ( $j \in \{0, 1, 2\}$ ), whose value varies from line to line. We establish the following lemmas to upper bound the terms on the right-hand side of (D.8).

**Lemma D.2** (Upper Bound of PDE – IP). Let  $f$  be any continuous function such that  $\|f\|_\infty \leq 1$  and  $\text{Lip}(f) \leq 1$ . Under Assumptions 4.1 and 4.2, it holds for any  $f$  that

$$\sup_{t \in [0, T]} \left| \int f(\theta) d\rho_t(\theta) - m^{-1} \cdot \sum_{i=1}^m f(\bar{\theta}_i(t)) \right| \leq B \cdot \sqrt{\log(mT/\delta)/m}$$

with probability at least  $1 - \delta$ .

*Proof.* See §D.1.1 for a detailed proof. □

**Lemma D.3** (Upper Bound of IP – CTTD). Under Assumptions 4.1 and 4.2, it holds that

$$\sup_{t \in [0, T]} \|\bar{\theta}^{(m)}(t) - \tilde{\theta}^{(m)}(t)\|_{(m)} \leq B \cdot e^{BT} \cdot \sqrt{\log(m/\delta)/m}$$

with probability at least  $1 - \delta$ .

*Proof.* See §D.1.2 for a detailed proof. □

**Lemma D.4** (Upper Bound of CTTD – ETD). Under Assumptions 4.1 and 4.2, it holds that

$$\sup_{\substack{k \leq T/\epsilon \\ (k \in \mathbb{N})}} \|\tilde{\theta}^{(m)}(k\epsilon) - \check{\theta}^{(m)}(k)\|_{(m)} \leq B \cdot e^{BT} \cdot \epsilon.$$

*Proof.* See §D.1.3 for a detailed proof. □

**Lemma D.5** (Upper Bound of ETD – TD). Under Assumptions 4.1 and 4.2, it holds that

$$\sup_{\substack{k \leq T/\epsilon \\ (k \in \mathbb{N})}} \|\check{\theta}^{(m)}(k) - \theta^{(m)}(k)\|_{(m)} \leq B \cdot e^{BT} \cdot \sqrt{\epsilon \cdot (D + \log(m/\delta))}$$

with probability at least  $1 - \delta$

*Proof.* See §D.1.4 for a detailed proof.  $\square$

We are now ready to present the proof of Proposition D.1.

*Proof.* Plugging Lemmas D.2-D.5 into (D.7), we have that

$$\begin{aligned} & \sup_{\substack{k \leq T/\epsilon \\ (k \in \mathbb{N})}} \left| \int f(\theta) d\rho_{k\epsilon}(\theta) - \int f(\theta) d\hat{\rho}_k(\theta) \right| \\ & \leq B \cdot e^{BT} \cdot \left( \sqrt{\log(m/\delta)/m} + \sqrt{\epsilon \cdot (D + \log(m/\delta))} \right) \end{aligned}$$

with probability at least  $1 - \delta$ . Thus, we complete the proof of Proposition D.1.  $\square$

## D.1 Proofs of Lemmas D.2-D.5

In this section, we present the proofs of Lemmas D.2-D.5, which are based on [6, 53, 54]. We include the required technical lemmas in §D.3. Recall that  $B > 0$  is a constant that depends on  $\alpha, \eta, \gamma, B_r$ , and  $B_j$  ( $j \in \{0, 1, 2\}$ ), whose value varies from line to line.

### D.1.1 Proof of Lemma D.2

*Proof.* For the IP dynamics in (D.6), it holds that  $\bar{\theta}_i(t) \sim \rho_t$  ( $i \in [m]$ ) (Proposition 8.1.8 in [5]). Furthermore, since the randomness of  $\bar{\theta}_i(t)$  comes from  $\theta_i(0)$  while  $\theta_i(0)$  ( $i \in [m]$ ) are independent, we have that  $\bar{\theta}_i(t) \stackrel{\text{i.i.d.}}{\sim} \rho_t$  ( $i \in [m]$ ). Thus, we have that

$$\mathbb{E}_{\rho_t} \left[ m^{-1} \cdot \sum_{i=1}^m f(\bar{\theta}_i(t)) \right] = \int f(\theta) d\rho_t(\theta).$$

Let  $\theta^{1,(m)} = \{\theta_1, \dots, \theta_i^1, \dots, \theta_m\}$  and  $\theta^{2,(m)} = \{\theta_1, \dots, \theta_i^2, \dots, \theta_m\}$  be two sets that only differ in the  $i$ -th element. Then, by the condition of Lemma D.2 that  $\|f\|_\infty \leq 1$ , we have that

$$\left| m^{-1} \cdot \sum_{j=1}^m f(\theta_j^1) - m^{-1} \cdot \sum_{j=1}^m f(\theta_j^2) \right| = m^{-1} \cdot |f(\theta_i^1) - f(\theta_i^2)| \leq 2/m.$$

Applying McDiarmid's inequality [70], we have for a fixed  $t \in [0, T]$  that

$$\mathbb{P} \left( \left| m^{-1} \cdot \sum_{i=1}^m f(\bar{\theta}_i(t)) - \int f(\theta) d\rho_t(\theta) \right| \geq p \right) \leq \exp(-mp^2/4). \quad (\text{D.9})$$

Moreover, we have for any  $s, t \in [0, T]$  that

$$\begin{aligned} & \left| m^{-1} \cdot \sum_{i=1}^m f(\bar{\theta}_i(t)) - \int f(\theta) d\rho_t(\theta) \right| - \left| m^{-1} \cdot \sum_{i=1}^m f(\bar{\theta}_i(s)) - \int f(\theta) d\rho_s(\theta) \right| \\ & \leq \left| m^{-1} \cdot \sum_{i=1}^m f(\bar{\theta}_i(t)) - m^{-1} \cdot \sum_{i=1}^m f(\bar{\theta}_i(s)) \right| + \left| \int f(\theta) d\rho_t(\theta) - \int f(\theta) d\rho_s(\theta) \right| \\ & \leq \|\bar{\theta}^{(m)}(t) - \bar{\theta}^{(m)}(s)\|_{(m)} + \mathcal{W}_1(\rho_t, \rho_s) \\ & \leq \|\bar{\theta}^{(m)}(t) - \bar{\theta}^{(m)}(s)\|_{(m)} + \mathcal{W}_2(\rho_t, \rho_s), \end{aligned}$$

where the second inequality follows from the fact that  $\text{Lip}(f) \leq 1$  and (C.9). Applying (D.38) and (D.40) of Lemma D.8, we have for any  $s, t \in [0, T]$  that

$$\left| m^{-1} \cdot \sum_{i=1}^m f(\bar{\theta}_i(t)) - \int f(\theta) d\rho_t(\theta) \right| - \left| m^{-1} \cdot \sum_{i=1}^m f(\bar{\theta}_i(s)) - \int f(\theta) d\rho_s(\theta) \right| \leq B \cdot |t - s|.$$

Applying the union bound to (D.9) for  $t \in \iota \cdot \{0, 1, \dots, \lfloor T/\iota \rfloor\}$ , we have that

$$\mathbb{P}\left(\sup_{t \in [0, T]} \left| m^{-1} \cdot \sum_{i=1}^m f(\bar{\theta}_i(t)) - \int f(\theta) d\rho_t(\theta) \right| \geq p + B \cdot \iota \right) \leq (T/\iota + 1) \cdot \exp(-mp^2/4).$$

Setting  $\iota = m^{-1/2}$  and  $p = B \cdot \sqrt{\log(mT/\delta)/m}$ , we have that

$$\sup_{t \in [0, T]} \left| m^{-1} \cdot \sum_{i=1}^m f(\bar{\theta}_i(t)) - \int f(\theta) d\rho_t(\theta) \right| \leq B \cdot \sqrt{\log(mT/\delta)/m}$$

with probability at least  $1 - \delta$ . Thus, we complete the proof of Lemma D.2.  $\square$

### D.1.2 Proof of Lemma D.3

*Proof.* Recall that  $g$  and  $\hat{g}$  are defined in (3.5) and (D.1), respectively, that is,

$$\begin{aligned} g(\theta; \rho) &= -\alpha \cdot \mathbb{E}_{\mathcal{D}} \left[ (Q(x; \rho) - r - \gamma \cdot Q(x'; \rho)) \cdot \nabla_{\theta} \sigma(x; \theta) \right], \\ \hat{g}(\theta; \theta^{(m)}) &= -\alpha \cdot \mathbb{E}_{\mathcal{D}} \left[ (\hat{Q}(x; \theta^{(m)}) - r - \gamma \cdot \hat{Q}(x'; \theta^{(m)})) \cdot \nabla_{\theta} \sigma(x; \theta) \right]. \end{aligned}$$

Following from the definition of  $\tilde{\theta}_i(t)$  and  $\bar{\theta}_i(t)$  in (D.5) and (D.6), respectively, we have for any  $i \in [m]$  and  $t \in [0, T]$  that

$$\begin{aligned} & \|\bar{\theta}_i(t) - \tilde{\theta}_i(t)\| \\ & \leq \int_0^t \left\| \frac{d\tilde{\theta}_i(s)}{ds} - \frac{d\bar{\theta}_i(s)}{ds} \right\| ds \\ & = \eta \cdot \int_0^t \left\| \hat{g}(\tilde{\theta}_i(s); \tilde{\theta}^{(m)}(s)) - g(\bar{\theta}_i(s); \rho_s) \right\| ds \\ & \leq \eta \cdot \int_0^t \left\| \hat{g}(\tilde{\theta}_i(s); \tilde{\theta}^{(m)}(s)) - \hat{g}(\bar{\theta}_i(s); \bar{\theta}^{(m)}(s)) \right\| ds + \eta \cdot \int_0^t \left\| \hat{g}(\bar{\theta}_i(s); \bar{\theta}^{(m)}(s)) - g(\bar{\theta}_i(s); \rho_s) \right\| ds \\ & \leq B \cdot \int_0^t \left\| \tilde{\theta}^{(m)}(s) - \bar{\theta}^{(m)}(s) \right\|_{(m)} ds + \eta \cdot \int_0^t \left\| \hat{g}(\bar{\theta}_i(s); \bar{\theta}^{(m)}(s)) - g(\bar{\theta}_i(s); \rho_s) \right\| ds, \quad (\text{D.10}) \end{aligned}$$

where the last inequality follows from (D.35) of Lemma D.7. We now upper bound the second term on the right-hand side of (D.10). Following from the definition of  $\hat{Q}$ ,  $Q$ , and  $\hat{g}$  in (3.1), (3.2), and (D.1), respectively, we have for any  $s \in [0, T]$  and  $i \in [m]$  that

$$\left\| \hat{g}(\bar{\theta}_i(s); \bar{\theta}^{(m)}(s)) - g(\bar{\theta}_i(s); \rho_s) \right\| = \alpha^2 \cdot \left\| m^{-1} \cdot \sum_{j=1}^m Z_i^j(s) \right\|, \quad (\text{D.11})$$

where

$$Z_i^j(s) = \mathbb{E}_{\mathcal{D}} \left[ \left( \sigma(x; \bar{\theta}_j(s)) - \int \sigma(x; \theta) d\rho_s(\theta) - \gamma \cdot \sigma(x'; \bar{\theta}_j(s)) + \gamma \cdot \int \sigma(x'; \theta) d\rho_s(\theta) \right) \cdot \nabla_{\theta} \sigma(x; \bar{\theta}_i(s)) \right].$$

Following from Assumptions 4.1 and 4.2, we have that  $\|Z_i^j(s)\| \leq B$ . When  $i \neq j$ , following from the fact that  $\bar{\theta}_i(s) \stackrel{\text{i.i.d.}}{\sim} \rho_s$  ( $i \in [m]$ ), it holds that  $\mathbb{E}[Z_i^j(s) | \bar{\theta}_i(s)] = 0$ . Following from Lemma D.9, we have for fixed  $s \in [0, T]$  and  $i \in [m]$  that

$$\begin{aligned} \mathbb{P}\left(\left\| m^{-1} \cdot \sum_{j \neq i} Z_i^j(s) \right\| \geq B \cdot (m^{-1/2} + p)\right) &= \mathbb{E}\left[\mathbb{P}\left(\left\| m^{-1} \cdot \sum_{j \neq i} Z_i^j(s) \right\| \geq B \cdot (m^{-1/2} + p) \mid \bar{\theta}_i(s)\right)\right] \\ &\leq \exp(-mp^2). \quad (\text{D.12}) \end{aligned}$$

By (C.9), we have that

$$\sup_{x \in \mathcal{X}} \left| \int \sigma(x; \theta) d\rho_s(\theta) - \int \sigma(x; \theta) d\rho_t(\theta) \right| \leq B \cdot \mathcal{W}_1(\rho_s, \rho_t) \leq B \cdot \mathcal{W}_2(\rho_s, \rho_t) \leq B \cdot |s - t|,$$

where the last inequality follows from (D.40) of Lemma D.8. Thus, following from Assumptions 4.1 and 4.2, Lemma D.8, and the fact that  $\text{Lip}(fg) \leq \|f\|_\infty \cdot \text{Lip}(g) + \|g\|_\infty \cdot \text{Lip}(f)$  for any functions  $f$  and  $g$ , we have for any  $s, t \in [0, T]$  that

$$\left\| m^{-1} \cdot \sum_{j \neq i} Z_i^j(s) \right\| - \left\| m^{-1} \cdot \sum_{j \neq i} Z_i^j(t) \right\| \leq B \cdot |t - s|.$$

Applying the union bound to (D.12) for  $i \in [m]$  and  $t \in \iota \cdot \{0, 1, \dots, \lfloor T/\iota \rfloor\}$ , we have that

$$\mathbb{P} \left( \sup_{\substack{i \in [m], \\ s \in [0, T]}} \left\| m^{-1} \cdot \sum_{j \neq i} Z_i^j(s) \right\| \geq B \cdot (m^{-1/2} + p) + B\iota \right) \leq m \cdot (T/\iota + 1) \cdot \exp(-mp^2).$$

Setting  $\iota = m^{-1/2}$  and  $p = B \cdot \sqrt{\log(mT/\delta)/m}$ , we have that

$$\sup_{\substack{i \in [m], \\ s \in [0, T]}} \left\| m^{-1} \cdot \sum_{j \neq i} Z_i^j(s) \right\| \leq B \cdot \sqrt{\log(mT/\delta)/m} \quad (\text{D.13})$$

with probability at least  $1 - \delta$ . When  $i = j$ , it holds that  $\|m^{-1} \cdot Z_i^i(s)\| \leq B/m$  in (D.11), which follows from Assumptions 4.1 and 4.2. Thus, plugging (D.13) into (D.11), we have that

$$\begin{aligned} \sup_{\substack{i \in [m], \\ s \in [0, T]}} \left\| \hat{g}(\bar{\theta}_i(s); \bar{\theta}^{(m)}(s)) - g(\bar{\theta}_i(s); \rho_s) \right\| &\leq \sup_{\substack{i \in [m], \\ s \in [0, T]}} \alpha^2 \cdot \left( \left\| m^{-1} \cdot Z_i^i(s) \right\| + \left\| m^{-1} \cdot \sum_{j \neq i} Z_i^j(s) \right\| \right) \\ &\leq B \cdot \sqrt{\log(mT/\delta)/m} \end{aligned} \quad (\text{D.14})$$

with probability at least  $1 - \delta$ .

Conditioning on the event in (D.14), we obtain from (D.10) that

$$\left\| \tilde{\theta}^{(m)}(t) - \bar{\theta}^{(m)}(t) \right\|_{(m)} \leq B \cdot \int_0^t \left\| \tilde{\theta}^{(m)}(s) - \bar{\theta}^{(m)}(s) \right\|_{(m)} ds + BT \cdot \sqrt{\log(mT/\delta)/m}$$

for any  $t \in [0, T]$ . Following from Gronwall's Lemma [41], we have that

$$\begin{aligned} \left\| \tilde{\theta}^{(m)}(t) - \bar{\theta}^{(m)}(t) \right\|_{(m)} &\leq B \cdot e^{Bt} \cdot BT \cdot \sqrt{\log(mT/\delta)/m} \\ &\leq B \cdot e^{BT} \cdot \sqrt{\log(m/\delta)/m}, \quad \forall t \in [0, T] \end{aligned}$$

with probability at least  $1 - \delta$ . Here the last inequality holds since we allow the value of  $B$  to vary from line to line. Thus, we complete the proof of Lemma D.3  $\square$

### D.1.3 Proof of Lemma D.4

*Proof.* By the definition of  $\hat{g}$ ,  $\check{\theta}_i(t)$ , and  $\tilde{\theta}_i(t)$  in (D.1), (D.4), and (D.5), respectively, it holds that

$$\begin{aligned} \left\| \tilde{\theta}_i(k\epsilon) - \check{\theta}_i(k) \right\| &\leq \eta \cdot \int_0^{k\epsilon} \left\| \hat{g}(\tilde{\theta}_i(s); \tilde{\theta}^{(m)}(s)) - \hat{g}(\check{\theta}_i(\lfloor s/\epsilon \rfloor); \check{\theta}^{(m)}(\lfloor s/\epsilon \rfloor)) \right\| ds \\ &\leq \eta \cdot \int_0^{k\epsilon} \left\| \hat{g}(\tilde{\theta}_i(s); \tilde{\theta}^{(m)}(s)) - \hat{g}(\tilde{\theta}_i(\lfloor s/\epsilon \rfloor \cdot \epsilon); \tilde{\theta}^{(m)}(\lfloor s/\epsilon \rfloor \cdot \epsilon)) \right\| ds \\ &\quad + \eta \cdot \sum_{\ell=0}^{k-1} \left\| \hat{g}(\tilde{\theta}_i(\ell\epsilon); \tilde{\theta}^{(m)}(\ell\epsilon)) - \hat{g}(\check{\theta}_i(\ell); \check{\theta}^{(m)}(\ell)) \right\| \\ &\leq B \cdot k \cdot \epsilon^2 + B \cdot \sum_{\ell=0}^{k-1} \left\| \tilde{\theta}^{(m)}(\ell\epsilon) - \check{\theta}^{(m)}(\ell) \right\|_{(m)}, \end{aligned}$$

where the last inequality follows from (D.35) of Lemma D.7 and (D.39) of Lemma D.8. Following from the definition of  $\|\cdot\|_{(m)}$  in (D.8), it holds for any  $k \leq T/\epsilon$  ( $k \in \mathbb{N}$ ) that

$$\left\| \tilde{\theta}^{(m)}(k\epsilon) - \check{\theta}^{(m)}(k) \right\|_{(m)} \leq B \cdot T \cdot \epsilon + B \cdot \sum_{\ell=0}^{k-1} \left\| \tilde{\theta}^{(m)}(\ell\epsilon) - \check{\theta}^{(m)}(\ell) \right\|_{(m)}.$$

Following from the discrete Gronwall's lemma [41], we have that

$$\sup_{\substack{k \leq T/\epsilon \\ (k \in \mathbb{N})}} \|\tilde{\theta}^{(m)}(k\epsilon) - \check{\theta}^{(m)}(k)\|_{(m)} \leq B^2 \cdot T \cdot \epsilon \cdot e^{BT} \leq B \cdot e^{BT} \cdot \epsilon,$$

where the last inequality holds since we allow the value of  $B$  to vary from line to line. Thus, we complete the proof of Lemma D.4.  $\square$

#### D.1.4 Proof of Lemma D.5

*Proof.* Let  $\mathcal{G}_k = \sigma(\theta^{(m)}(0), z_0, \dots, z_k)$  be the  $\sigma$ -algebra generated by  $\theta^{(m)}(0)$  and  $z_\ell = (x_\ell, r_\ell, x'_\ell)$  ( $\ell \leq k$ ). Recall that  $\hat{g}$  and  $\hat{G}_k$  are defined in (D.1) and (D.2), respectively. We have for any  $i \in [m]$  and  $k \in \mathbb{N}_+$  that

$$\mathbb{E}[\hat{G}_k(\theta_i(k); \theta^{(m)}(k)) \mid \mathcal{G}_{k-1}] = \hat{g}(\theta_i(k); \theta^{(m)}(k)).$$

Recall that  $\theta^{(m)}(k)$  and  $\check{\theta}^{(m)}(k)$  are the TD and ETD dynamics defined in (D.3) and (D.4), respectively. Thus, we have for any  $i \in [m]$  and  $k \in \mathbb{N}_+$  that

$$\begin{aligned} \|\check{\theta}_i(k) - \theta_i(k)\| &= \eta\epsilon \cdot \left\| \sum_{\ell=0}^{k-1} \hat{G}_\ell(\theta_i(\ell); \theta^{(m)}(\ell)) - \sum_{\ell=0}^{k-1} \hat{g}(\check{\theta}_i(\ell); \check{\theta}^{(m)}(\ell)) \right\| \\ &\leq \eta\epsilon \cdot \left\| \sum_{\ell=0}^{k-1} X_i(\ell) \right\| + \eta\epsilon \cdot \sum_{\ell=0}^{k-1} \left\| \hat{g}(\check{\theta}_i(\ell); \check{\theta}^{(m)}(\ell)) - \hat{g}(\theta_i(\ell); \theta^{(m)}(\ell)) \right\| \\ &\leq \eta\epsilon \cdot \|A_i(k)\| + B\epsilon \cdot \sum_{\ell=0}^{k-1} \|\check{\theta}^{(m)}(\ell) - \theta^{(m)}(\ell)\|_{(m)}, \end{aligned} \quad (\text{D.15})$$

where the last inequality follows from (D.35) of Lemma D.7, and  $X_i(\ell)$  and  $A_i(k)$  are defined as

$$\begin{aligned} X_i(0) &= 0, \\ X_i(\ell) &= \hat{G}_\ell(\theta_i(\ell); \theta^{(m)}(\ell)) - \mathbb{E}[\hat{G}_\ell(\theta_i(\ell); \theta^{(m)}(\ell)) \mid \mathcal{G}_{\ell-1}] \quad \forall \ell \geq 1, \\ A_i(k) &= \sum_{\ell=0}^{k-1} X_i(\ell). \end{aligned}$$

Following from (D.32) of Lemma D.7, we have that  $\|X_i(\ell)\| \leq B$ . Thus, the stochastic process  $\{A_i(k)\}_{k \in \mathbb{N}_+}$  is a martingale with  $\|A_i(k) - A_i(k-1)\| \leq B$ . Applying Lemma D.10, we have that

$$\mathbb{P}\left(\max_{\substack{k \leq T/\epsilon \\ (k \in \mathbb{N}_+)}} \|A_i(k)\| \geq B \cdot \sqrt{T/\epsilon} \cdot (\sqrt{D} + p)\right) \leq \exp(-p^2). \quad (\text{D.16})$$

Applying the union bound to (D.16) for  $i \in [m]$ , we have that

$$\mathbb{P}\left(\max_{\substack{i \in [m], \\ k \leq T/\epsilon \ (k \in \mathbb{N}_+)}} \|A_i(k)\| \geq B \cdot \sqrt{T/\epsilon} \cdot (\sqrt{D} + p)\right) \leq m \cdot \exp(-p^2).$$

By setting  $p = \sqrt{\log(m/\delta)}$ , we have that

$$\|A_i(k)\| \leq B \cdot \sqrt{T/\epsilon} \cdot (\sqrt{D} + \sqrt{\log(m/\delta)}), \quad \forall i \in [m], k \leq T/\epsilon \ (k \in \mathbb{N}_+) \quad (\text{D.17})$$

with probability at least  $1 - \delta$ . By (D.15) and (D.17), we have that

$$\begin{aligned} \|\check{\theta}^{(m)}(k) - \theta^{(m)}(k)\|_{(m)} &\leq B \cdot \sqrt{T\epsilon} \cdot (\sqrt{D} + \sqrt{\log(m/\delta)}) + B\epsilon \cdot \sum_{\ell=0}^{k-1} \|\check{\theta}^{(m)}(\ell) - \theta^{(m)}(\ell)\|_{(m)}, \quad \forall k \leq T/\epsilon \ (k \in \mathbb{N}) \end{aligned}$$

with probability at least  $1 - \delta$ . Applying the discrete Gronwall's Lemma [41], we have that

$$\begin{aligned} \|\check{\theta}^{(m)}(k) - \theta^{(m)}(k)\|_{(m)} &\leq B \cdot e^{BT} \cdot B \cdot \sqrt{T\epsilon} \cdot (\sqrt{D} + \sqrt{\log(m/\delta)}) \\ &\leq B \cdot e^{BT} \cdot \sqrt{\epsilon \cdot (D + \log(m/\delta))}, \quad \forall k \leq T/\epsilon \ (k \in \mathbb{N}) \end{aligned}$$

with probability at least  $1 - \delta$ . Here the last inequality holds since we allow the value of  $B$  to vary from line to line. Thus, we complete the proof of Lemma D.5.  $\square$

## D.2 Proof of Corollary 4.4

The proof of Corollary 4.4 follows from Theorem 4.3 and the following lemma, which characterizes the error of approximating the TD dynamics  $\theta^{(m)}(k)$  in (3.3) using the PDE solution  $\rho_t$  in (3.4).

**Lemma D.6.** Let  $B$  be a constant that depends on  $\alpha, \eta, \gamma, B_0, B_1$ , and  $B_2$ . Under Assumptions 4.1 and 4.2, it holds for any  $k \leq T/\epsilon$  ( $k \in \mathbb{N}$ ) that

$$\begin{aligned} & \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( \widehat{Q}(x; \theta^{(m)}(k)) - Q^*(x) \right)^2 \right] \\ & \leq \mathbb{E}_{x \sim \mathcal{D}} \left[ \left( Q(x; \rho_{k\epsilon}) - Q^*(x) \right)^2 \right] + B \cdot e^{BT} \cdot \left( \sqrt{m^{-1} \cdot \log(m/\delta)} + \sqrt{\epsilon \cdot (D + \log(m/\delta))} \right) \end{aligned}$$

with probability at least  $1 - \delta$ .

*Proof.* Recall that  $\widehat{Q}$  and  $Q(\cdot; \rho)$  are defined in (3.1) and (3.2), respectively. For notational simplicity, we denote the optimality gaps for  $\theta^{(m)} = \{\theta_i\}_{i=1}^m$  and  $\rho \in \mathcal{P}_2(\mathbb{R}^D)$  by

$$L(\theta^{(m)}) = \mathbb{E}_{\mathcal{D}} \left[ \left( \widehat{Q}(x; \theta^{(m)}) - Q^*(x) \right)^2 \right], \quad (\text{D.18})$$

$$\bar{L}(\rho) = \mathbb{E}_{\mathcal{D}} \left[ \left( Q(x; \rho) - Q^*(x) \right)^2 \right]. \quad (\text{D.19})$$

Recall that  $\theta^{(m)}(k)$ ,  $\bar{\theta}^{(m)}(k\epsilon)$ , and  $\rho_t$  are the TD dynamics, the IP dynamics, and the PDE solution defined in (D.3), (D.6), and (3.4), respectively. It holds for any  $k \in \mathbb{N}$  that

$$\left| L(\theta^{(m)}(k)) - \bar{L}(\rho_{k\epsilon}) \right| \leq \underbrace{\left| L(\theta^{(m)}(k)) - L(\bar{\theta}^{(m)}(k\epsilon)) \right|}_{(\text{i})} + \underbrace{\left| L(\bar{\theta}^{(m)}(k\epsilon)) - \bar{L}(\rho_{k\epsilon}) \right|}_{(\text{ii})}. \quad (\text{D.20})$$

In what follows, we upper bound the two terms on the right-hand side of (D.20).

**Upper bounding term (i) of (D.20).** Following from the definition of  $L$  in (D.18), it holds for any  $k \in \mathbb{N}$  that

$$\begin{aligned} & \left| L(\theta^{(m)}(k)) - L(\bar{\theta}^{(m)}(k\epsilon)) \right| \\ & = \left| \mathbb{E}_{\mathcal{D}} \left[ \left( \widehat{Q}(x; \theta^{(m)}(k)) + \widehat{Q}(x; \bar{\theta}_i(k\epsilon)) - 2Q^*(x) \right) \cdot \left( \widehat{Q}(x; \theta^{(m)}(k)) - \widehat{Q}(x; \bar{\theta}_i(k\epsilon)) \right) \right] \right|. \end{aligned} \quad (\text{D.21})$$

Following from (D.30), (D.31), and (D.36) of Lemma D.7, we have for any  $k \in \mathbb{N}$  that

$$\sup_{x \in \mathcal{X}} \left| \widehat{Q}(x; \theta^{(m)}(k)) + \widehat{Q}(x; \bar{\theta}_i(k\epsilon)) - 2Q^*(x) \right| \leq B, \quad (\text{D.22})$$

$$\sup_{x \in \mathcal{X}} \left| \widehat{Q}(x; \theta^{(m)}(k)) - \widehat{Q}(x; \bar{\theta}_i(k\epsilon)) \right| \leq B \cdot \|\theta^{(m)}(k) - \bar{\theta}^{(m)}(k\epsilon)\|_{(m)}. \quad (\text{D.23})$$

Thus, we have that

$$\begin{aligned} & \left| L(\theta^{(m)}(k)) - L(\bar{\theta}^{(m)}(k\epsilon)) \right| \\ & \leq B \cdot \|\theta^{(m)}(k) - \bar{\theta}^{(m)}(k\epsilon)\|_{(m)} \\ & \leq B \cdot e^{BT} \cdot \left( \sqrt{\log(m/\delta)/m} + \sqrt{\epsilon \cdot (D + \log(m/\delta))} \right), \quad \forall k \leq T/\epsilon \ (k \in \mathbb{N}) \end{aligned} \quad (\text{D.24})$$

with probability at least  $1 - \delta$ . Here the last inequality follows from Lemmas D.3-D.5.

**Upper bounding term (ii) of (D.20).** Let  $t = k\epsilon$ . It holds for any  $t \in [0, T]$  that

$$\left| L(\bar{\theta}^{(m)}(t)) - \bar{L}(\rho_t) \right| \leq \left| L(\bar{\theta}^{(m)}(t)) - \mathbb{E}_{\rho_t} \left[ L(\bar{\theta}^{(m)}(t)) \right] \right| + \left| \mathbb{E}_{\rho_t} \left[ L(\bar{\theta}^{(m)}(t)) \right] - \bar{L}(\rho_t) \right|, \quad (\text{D.25})$$

where the expectation is with respect to  $\bar{\theta}_i(t) \stackrel{\text{i.i.d.}}{\sim} \rho_t$  ( $i \in [m]$ ). For the second term on the right-hand side of (D.25), following from the fact that  $\mathbb{E}_{\rho_t}[\hat{Q}(x; \bar{\theta}^{(m)}(t))] = Q(x; \rho_t)$  for any  $x \in \mathcal{X}$ , we have that

$$\begin{aligned} \left| \mathbb{E}_{\rho_t} [L(\bar{\theta}^{(m)}(t))] - \bar{L}(\rho_t) \right| &= \left| \int \mathbb{E}_{\rho_t} [\hat{Q}(x; \bar{\theta}^{(m)}(t))^2 - Q(x; \rho_t)^2] d\mathcal{D}(x) \right| \\ &= \left| \int \text{Var}_{\rho_t} [\hat{Q}(x; \bar{\theta}^{(m)}(t))] d\mathcal{D}(x) \right| \\ &\leq B/m, \end{aligned} \quad (\text{D.26})$$

where the inequality follows from the fact that  $\|\sigma\| \leq B$  in Assumption 4.2 and the independence of  $\bar{\theta}_i(t)$  ( $i \in [m]$ ). Let  $\theta^{1,(m)} = \{\theta_1, \dots, \theta_i^1, \dots, \theta_m\}$  and  $\theta^{2,(m)} = \{\theta_1, \dots, \theta_i^2, \dots, \theta_m\}$  be two sets that only differ in the  $i$ -th element. It holds that

$$|L(\theta^{1,(m)}) - L(\theta^{2,(m)})| \leq B \cdot m^{-1} \cdot \mathbb{E}_{\mathcal{D}} [|\sigma(x; \theta_i^1) - \sigma(x; \theta_i^2)|] \leq B/m,$$

where the first inequality follows from (D.21) and (D.22) and the second inequality follows from Assumption 4.2. Applying McDiarmid's inequality [70], we have for a fixed  $t \in [0, T]$  that

$$\mathbb{P} \left( \left| L(\bar{\theta}^{(m)}(t)) - \mathbb{E}_{\rho_t} [L(\bar{\theta}^{(m)}(t))] \right| \geq p \right) \leq \exp(-mp^2/B). \quad (\text{D.27})$$

It holds for any  $s, t \in [0, T]$  that

$$\begin{aligned} &\left| \left| L(\bar{\theta}^{(m)}(t)) - \mathbb{E}_{\rho_t} [L(\bar{\theta}^{(m)}(t))] \right| - \left| L(\bar{\theta}^{(m)}(s)) - \mathbb{E}_{\rho_s} [L(\bar{\theta}^{(m)}(s))] \right| \right| \\ &\leq B \cdot \|\bar{\theta}^{(m)}(t) - \bar{\theta}^{(m)}(s)\|_{(m)} \leq B \cdot |t - s|, \end{aligned}$$

where the first inequality follows from (D.21), (D.22), and (D.23) and the second inequality follows from (D.38) of Lemma D.8. Applying the union bound to (D.27) for  $t \in \iota \cdot \{0, 1, \dots, \lfloor T/\iota \rfloor\}$ , we have that

$$\mathbb{P} \left( \sup_{t \in [0, T]} \left| L(\bar{\theta}^{(m)}(t)) - \mathbb{E}_{\rho_t} [L(\bar{\theta}^{(m)}(t))] \right| \geq p + B\iota \right) \leq (T/\iota + 1) \cdot \exp(-mp^2/B),$$

Setting  $\iota = m^{-1/2}$  and  $p = B \cdot \sqrt{\log(mT\delta)/m}$ , we have that

$$\sup_{t \in [0, T]} \left| L(\bar{\theta}^{(m)}(t)) - \mathbb{E}_{\rho_t} [L(\bar{\theta}^{(m)}(t))] \right| \leq B \cdot \sqrt{\log(mT\delta)/m} \quad (\text{D.28})$$

with probability at least  $1 - \delta$ . Plugging (D.26) and (D.28) into (D.25), noting that  $t = k\epsilon$ , we have that

$$\left| L(\bar{\theta}^{(m)}(k\epsilon)) - \bar{L}(\rho_{k\epsilon}) \right| \leq B \cdot \sqrt{\log(mT\delta)/m}, \quad \forall k \leq T/\epsilon \ (k \in \mathbb{N}) \quad (\text{D.29})$$

with probability at least  $1 - \delta$ .

Plugging (D.24) and (D.29) into (D.20), we have that

$$\left| L(\theta^{(m)}(k)) - \bar{L}(\rho_{k\epsilon}) \right| \leq B \cdot e^{BT} \cdot \left( \sqrt{\log(m/\delta)/m} + \sqrt{\epsilon \cdot (D + \log(m/\delta))} \right), \quad \forall k \leq T/\epsilon \ (k \in \mathbb{N})$$

with probability at least  $1 - \delta$ . Thus, we complete the proof of Lemma D.6.  $\square$

### D.3 Technical Lemmas for §D

In what follows, we present the technical lemmas used in §D. Recall that  $\hat{Q}$ ,  $\hat{g}$ , and  $\hat{G}_k$  are defined in (3.1), (D.1), and (D.2), respectively. Let  $B > 0$  be a constant depending on  $\alpha, \eta, \gamma, B_r$ , and  $B_j$  ( $j \in \{0, 1, 2\}$ ), whose value varies from line to line.



**Lemma D.7.** Under Assumptions 4.1 and 4.2, it holds for  $\theta^{(m)} = \{\theta_i\}_{i=1}^m$  and  $\underline{\theta}^{(m)} = \{\underline{\theta}_i\}_{i=1}^m$  that

$$\sup_{x \in \mathcal{X}} |\widehat{Q}(x; \theta^{(m)})| \leq B, \quad (\text{D.30})$$

$$\sup_{x \in \mathcal{X}} |\widehat{Q}(x; \theta^{(m)}) - \widehat{Q}(x; \underline{\theta}^{(m)})| \leq B \cdot \|\theta^{(m)} - \underline{\theta}^{(m)}\|_{(m)}, \quad (\text{D.31})$$

$$\|\widehat{G}_k(\theta_i; \theta^{(m)})\| \leq B, \quad (\text{D.32})$$

$$\|\widehat{G}_k(\theta_i; \theta^{(m)}) - \widehat{G}_k(\underline{\theta}_i; \underline{\theta}^{(m)})\| \leq B \cdot \|\theta^{(m)} - \underline{\theta}^{(m)}\|_{(m)}, \quad \forall k \in \mathbb{N}, \quad (\text{D.33})$$

$$\|\widehat{g}(\theta_i; \theta^{(m)})\| \leq B, \quad (\text{D.34})$$

$$\|\widehat{g}(\theta_i; \theta^{(m)}) - \widehat{g}(\underline{\theta}_i; \underline{\theta}^{(m)})\| \leq B \cdot \|\theta^{(m)} - \underline{\theta}^{(m)}\|_{(m)}. \quad (\text{D.35})$$

Meanwhile, for any  $Q \in \mathcal{F}$ , it holds that

$$\sup_{x \in \mathcal{X}} \|Q(x)\| \leq B. \quad (\text{D.36})$$

For any  $\rho \in \mathcal{P}_2(\mathbb{R}^D)$ , it holds that

$$\|g(\theta; \rho)\| \leq B. \quad (\text{D.37})$$

*Proof.* For (D.30) and (D.31) of Lemma D.7, following from Assumptions 4.1 and 4.2 and the definition of  $\widehat{Q}$  in (3.1), we have for any  $x \in \mathcal{X}$ ,  $\theta^{(m)}$ , and  $\underline{\theta}^{(m)}$  that

$$\begin{aligned} |\widehat{Q}(x; \theta^{(m)})| &\leq \alpha \cdot m^{-1} \sum_{i=1}^m |\sigma(x; \theta_i)| \leq B, \\ |\widehat{Q}(x; \theta^{(m)}) - \widehat{Q}(x; \underline{\theta}^{(m)})| &\leq \alpha \cdot m^{-1} \sum_{i=1}^m |\sigma(x; \theta_i) - \sigma(x; \underline{\theta}_i)| \leq B \cdot \|\theta^{(m)} - \underline{\theta}^{(m)}\|_{(m)}. \end{aligned}$$

For (D.32) and (D.33) of Lemma D.7, following from the definition of  $\widehat{G}_k$  in (D.2), we have for any  $\theta^{(m)}$  and  $\underline{\theta}^{(m)}$  that

$$\begin{aligned} \|\widehat{G}_k(\theta_i; \theta^{(m)})\| &= \alpha \cdot |\widehat{Q}(x_k; \theta^{(m)}) - r_k - \gamma \cdot \widehat{Q}(x'_k; \theta^{(m)})| \cdot \|\nabla_{\theta} \sigma(x_k; \theta_i)\| \leq B, \\ \|\widehat{G}_k(\theta_i; \theta^{(m)}) - \widehat{G}_k(\underline{\theta}_i; \underline{\theta}^{(m)})\| &\leq \alpha \cdot \sup_{\theta^{(m)}} |\widehat{Q}(x_k; \theta^{(m)}) - r_k - \gamma \cdot \widehat{Q}(x'_k; \theta^{(m)})| \cdot \|\nabla_{\theta} \sigma(x_k; \theta_i) - \nabla_{\theta} \sigma(x_k; \underline{\theta}_i)\| \\ &\quad + \alpha \cdot |\widehat{Q}(x_k; \theta^{(m)}) - \gamma \cdot \widehat{Q}(x'_k; \theta^{(m)}) - \widehat{Q}(x_k; \underline{\theta}^{(m)}) + \gamma \cdot \widehat{Q}(x'_k; \theta^{(m)})| \cdot \sup_{\theta_i \in \mathbb{R}^D} \|\nabla_{\theta} \sigma(x_k; \theta_i)\| \\ &\leq B \cdot \|\theta^{(m)} - \underline{\theta}^{(m)}\|_{(m)}. \end{aligned}$$

The inequalities in (D.34) and (D.35) of Lemma D.7 for  $\widehat{g}$  follow from the fact that

$$\widehat{g}(\theta_i; \theta^{(m)}) = \mathbb{E}_{(x_k, r_k, x'_k) \sim \bar{\mathcal{D}}} [G_k(\theta_i; \theta^{(m)})].$$

The inequalities in (D.36) and (D.37) follow from the definition of  $\mathcal{F}$  and  $g$  in (4.3) and (3.5), respectively. Thus, we complete the proof of Lemma D.7.  $\square$

Recall that  $\rho_t$  is the PDE solution in (3.4) and  $\widetilde{\theta}^{(m)}(t)$  and  $\bar{\theta}^{(m)}(t)$  are the CTTD and IP dynamics defined in (D.5) and (D.6), respectively.

**Lemma D.8.** Under Assumptions 4.1 and 4.2, it holds for any  $s, t \in [0, T]$  that

$$\|\bar{\theta}^{(m)}(t) - \bar{\theta}^{(m)}(s)\|_{(m)} \leq B \cdot |t - s|, \quad (\text{D.38})$$

$$\|\widetilde{\theta}^{(m)}(t) - \widetilde{\theta}^{(m)}(s)\|_{(m)} \leq B \cdot |t - s|, \quad (\text{D.39})$$

$$\mathcal{W}_2(\rho_t, \rho_s) \leq B \cdot |t - s|. \quad (\text{D.40})$$

*Proof.* For (D.38) of Lemma D.8, by the definition of  $\bar{\theta}_i(t)$  in (D.6) and (D.37) of Lemma D.7, we have for any  $s, t \in [0, T]$  and  $i \in [m]$  that

$$\|\bar{\theta}_i(t) - \bar{\theta}_i(s)\| = \eta \cdot \int_s^t \|g(\bar{\theta}_i(\tau); \rho_\tau)\| d\tau \leq B \cdot |t - s|.$$

Similarly, for (D.39) of Lemma D.8, by the definition of  $\tilde{\theta}_i(t)$  in (D.5) and (D.34) of Lemma D.7, we have for any  $i \in [m]$  and  $s, t \in [0, T]$  that  $\|\tilde{\theta}_i(t) - \tilde{\theta}_i(s)\| \leq B \cdot |t - s|$ .

For (D.40) of Lemma D.8, following from the fact that  $\bar{\theta}_i(t) \stackrel{\text{i.i.d.}}{\sim} \rho_t$  ( $i \in [m]$ ) and the definition of  $\mathcal{W}_2$  in (2.4), it holds for any  $s, t \in [0, T]$  that

$$\mathcal{W}_2(\rho_t, \rho_s) \leq \mathbb{E} \left[ \|\bar{\theta}_i(t) - \bar{\theta}_i(s)\|^2 \right]^{1/2} \leq B \cdot |t - s|.$$

Thus, we complete the proof of Lemma D.8.  $\square$

**Lemma D.9** (Lemma 30 in [53]). Let  $\{X_i\}_{i=1}^m$  be i.i.d. random variables with  $\|X_i\| \leq \xi$  and  $\mathbb{E}[X_i] = 0$ . Then, it holds for any  $p > 0$  that

$$\mathbb{P} \left( \left\| m^{-1} \cdot \sum_{i=1}^m X_i \right\| \geq C\xi \cdot (m^{-1/2} + p) \right) \leq \exp(-mp^2),$$

where  $C > 0$  is an absolute constant.

**Lemma D.10** (Lemma A.3 in [6] and Lemma 31 in [53]). Let  $X_k \in \mathbb{R}^D$  ( $k \in \mathbb{N}$ ) be a martingale with respect to the filtration  $\mathcal{G}_k$  ( $k \geq 0$ ) with  $X_0 = 0$ . We assume for  $\xi > 0$  and any  $\lambda \in \mathbb{R}^D$  that

$$\mathbb{E} \left[ \exp(\langle \lambda, X_k - X_{k-1} \rangle) \mid \mathcal{G}_{k-1} \right] \leq \exp(\xi^2 \cdot \|\lambda\|^2 / 2).$$

Then, it holds that

$$\mathbb{P} \left( \max_{\substack{k \leq n \\ (k \in \mathbb{N})}} \|X_k\| \geq C\xi \cdot \sqrt{n} \cdot (\sqrt{D} + p) \right) \leq \exp(-p^2),$$

where  $C > 0$  is an absolute constant.

## E Auxiliary Lemmas

We use the definition of absolutely continuous curves in  $\mathcal{P}_2(\mathbb{R}^D)$  in [5].

**Definition E.1** (Absolutely Continuous Curve). Let  $\beta : [a, b] \rightarrow \mathcal{P}_2(\mathbb{R}^D)$  be a curve. Then,  $\beta$  is an absolutely continuous curve if there exists a square-integrable function  $f : [a, b] \rightarrow \mathbb{R}$  such that

$$\mathcal{W}_2(\beta_s, \beta_t) \leq \int_s^t f(\tau) d\tau$$

for any  $a \leq s < t \leq b$ .

Then, we have the following first variation formula.

**Lemma E.2** (First Variation Formula, Theorem 8.4.7 in [5]). Given  $\nu \in \mathcal{P}_2(\mathbb{R}^D)$  and an absolutely continuous curve  $\mu : [0, T] \rightarrow \mathcal{P}_2(\mathbb{R}^D)$ , let  $\beta : [0, 1] \rightarrow \mathcal{P}_2(\mathbb{R}^D)$  be the geodesic connecting  $\mu_t$  and  $\nu$ . It holds that

$$\frac{d}{dt} \frac{\mathcal{W}_2(\mu_t, \nu)^2}{2} = -\langle \mu'_t, \beta'_0 \rangle_{\mu_t},$$

where  $\mu'_t = \partial_t \mu_t$ ,  $\beta'_0 = \partial_t \beta_t|_{t=0}$ , and the inner product is defined in (2.5).

**Lemma E.3** (Talagrand's Inequality, Corollary 2.1 in [59]). Let  $\nu$  be  $N(0, \kappa \cdot I_D)$ . It holds for any  $\mu \in \mathcal{P}_2(\mathbb{R}^D)$  that

$$\mathcal{W}_2(\mu, \nu)^2 \leq 2D_{\text{KL}}(\mu \parallel \nu) / \kappa.$$

**Lemma E.4** (Eulerian Representation of Geodesics, Proposition 5.38 in [68]). Let  $\beta : [0, 1] \rightarrow \mathcal{P}_2(\mathbb{R}^D)$  be a geodesic and  $v$  be the corresponding vector field such that  $\partial_t \beta_t = -\text{div}(\beta_t \cdot v_t)$ . It holds that

$$\partial_t(\beta_t \cdot v_t) = -\text{div}(\beta_t \cdot v_t \otimes v_t).$$