

## 467 Appendix A Continuous RL: Formulation and Well-Posedness

### 468 A.1 Exploratory Stochastic-Control

469 For  $n, m$  positive integers, let  $b : \mathbb{R}^n \times \mathcal{A} \mapsto \mathbb{R}^n$  and  $\sigma : \mathbb{R}^n \times \mathcal{A} \mapsto \mathbb{R}^{n \times m}$  be given functions,  
 470 where  $\mathcal{A}$  is a compact action space. A classical stochastic control problem [15, 62] is to control  
 471 the state (or feature) dynamics governed by an Itô process, defined on a filtered probability space  
 472  $(\Omega, \mathcal{F}, \mathbb{P}; \{\mathcal{F}_s^B\}_{s \geq 0})$ , along with an  $\{\mathcal{F}_s^B\}$ -Brownian motion  $B = \{B_s, s \geq 0\}$ :

$$dX_s^a = b(X_s^a, a_s) ds + \sigma(X_s^a, a_s) dB_s, \quad s \geq t, \quad X_t = x, \quad (29)$$

473 where  $a_s$  is the agent's action (control) at time  $s$ . The goal of the stochastic control (discounted  
 474 objective over an infinite time horizon) is for any time-state pair  $(t, x)$  in (29), to find the optimal  
 475  $\{\mathcal{F}_s^B\}_{s \geq 0}$ -progressively measurable sequence of actions  $a = \{a_s, s \geq t\}$  (called the optimal policy)  
 476 that maximizes the expected total  $\beta$ -discounted reward:

$$\mathbb{E} \left[ \int_t^{+\infty} e^{-\beta(s-t)} r(X_s^a, a_s) ds \mid X_t^a = x \right], \quad (30)$$

477 where  $r : \mathbb{R}^n \times \mathcal{A} \mapsto \mathbb{R}$  is the running reward of the current state and action  $(X_s^a, a_s)$ , and  $\beta > 0$  is a  
 478 discount factor that measures the time-depreciation of the objective value (or the impatience level of  
 479 the agent). Note that the state process  $X^a = \{X_s^a, s \geq t\}$  depends on the starting (initial) time-state  
 480 pair  $(t, x)$ . For ease of notation, we denote by  $X^a$  instead of  $X^{t,x,a} = \{X_s^{t,x,a}, s \geq t\}$  the solution  
 481 to the SDE in (29) when there is no ambiguity.

482 Listed below are the standard assumptions to ensure the well-posedness of the stochastic control  
 483 problem in (29)-(30).

484 **Assumption 2.** *The following conditions are assumed throughout:*

- 485 (i)  $b, \sigma, r$  are all continuous functions in their respective arguments;  
 486 (ii)  $b, \sigma$  are uniformly Lipschitz continuous in  $x$ , i.e., there exists a constant  $C > 0$  such that for  
 487  $\varphi \in \{b, \sigma\}$ ,

$$\|\varphi(x, a) - \varphi(x', a)\|_2 \leq C \|x - x'\|_2, \quad \text{for all } a \in \mathcal{A}, \quad x, x' \in \mathbb{R}^n; \quad (31)$$

- 488 (iii)  $b, \sigma$  have linear growth in  $x$  and  $a$ , i.e., there exists a constant  $C > 0$  such that for  $\varphi \in \{b, \sigma\}$ ,

$$\|\varphi(x, a)\|_2 \leq C(1 + \|x\|_2 + \|a\|_2), \quad \text{for all } (x, a) \in \mathbb{R}^n \times \mathcal{A}; \quad (32)$$

- 489 (iv)  $r$  has polynomial growth in  $x$  and  $a$ , i.e., there exists a constant  $C > 0$  and  $\mu \geq 1$  such that

$$|r(x, a)| \leq C(1 + \|x\|_2^\mu + \|a\|_2^\mu) \quad \text{for all } (x, a) \in \mathbb{R}^n \times \mathcal{A}. \quad (33)$$

490 The key idea underlying *exploratory* stochastic control is to use a randomized policy (or relaxed  
 491 control), i.e., apply a probability distribution to the admissible action space. To do so, let's assume  
 492 the probability space is rich enough to support a uniform random variable  $Z$  that is independent  
 493 of the Brownian motion  $B = \{B_t\}$ . We then expand the original filtered probability space to  
 494  $(\Omega, \mathcal{F}, \mathbb{P}; \{\mathcal{F}_s\}_{s \geq 0})$ , where  $\mathcal{F}_s = \mathcal{F}_s^B \vee \sigma(Z)$  (i.e., augment  $\mathcal{F}_s^B$  with the sigma field generated by  
 495  $Z$ ).

496 Let  $\pi : \mathbb{R}^n \ni x \mapsto \pi(\cdot \mid x) \in \mathcal{P}(\mathcal{A})$  be a stationary feedback policy given the state at  $x$ , where  $\mathcal{P}(\mathcal{A})$   
 497 is a suitable collection of probability distributions (with density functions). At each time  $s$ , an action  
 498  $a_s$  is generated from the distribution  $\pi(\cdot \mid X_s^a)$ , i.e. the policy only depends on the current state.  
 499 In other words, we only consider stationary, or time-independent feedback control policies for the  
 500 stochastic control problem (29)-(30).

501 Given a stationary policy  $\pi \in \mathcal{P}(\mathcal{A})$ , an initial state  $x$ , and an  $\{\mathcal{F}_s\}$ -progressively measurable action  
 502 process  $a^\pi = \{a_s^\pi, s \geq 0\}$  generated from  $\pi$ , the state process  $X^\pi = \{X_s^\pi, s \geq 0\}$  follows:

$$dX_s^\pi = b(X_s^\pi, a_s^\pi) ds + \sigma(X_s^\pi, a_s^\pi) dB_s, \quad s \geq t, \quad X_0^\pi = x, \quad (34)$$

503 defined on  $(\Omega, \mathcal{F}, \mathbb{P}; \{\mathcal{F}_s\}_{s \geq 0})$ . It is easy to see that the dynamics in (34) define a time-homogeneous  
 504 Markov process, such that for each  $t \geq 0$  and  $x$ :

$$(X_s^\pi \mid X_0^\pi = x) \stackrel{d}{=} (X_{s+t}^\pi \mid X_t^\pi = x), \quad s \geq 0.$$

505 Consequently, the objective in (30) is independent of time  $t$ , and is equal to:

$$\mathbb{E} \left[ \int_0^{+\infty} e^{-\beta s} r(X_s^\pi, a_s^\pi) ds \mid X_0^\pi = x \right]. \quad (35)$$

506 Furthermore, following [58], we can add a regularizer to the objective function to encourage explo-  
507 ration (represented by the randomized policy), leading to

$$V(t, x; \pi) := \mathbb{E} \left[ \int_t^\infty e^{-\beta(s-t)} [r(X_s^\pi, a_s^\pi) + \gamma p(X_s^\pi, a_s^\pi, \pi(\cdot \mid X_s^\pi))] ds \mid X_t^\pi = x \right], \quad (36)$$

508 where  $p: \mathbb{R}^n \times \mathcal{A} \times \mathcal{P}(\mathcal{A}) \mapsto \mathbb{R}$  is the regularizer, and  $\gamma \geq 0$  is a weight parameter on exploration  
509 (also known as the ‘‘temperature’’ parameter). For instance, in [58],  $p$  is taken as the differential  
510 entropy,

$$p(x, a, \pi(\cdot)) := -\log \pi(a),$$

511 and hence, the ‘‘entropy’’ regularizer. The same argument as before justifies that  $V(t, x; \pi)$  is  
512 independent of time  $t$ . That is, for all  $t \geq 0$ ,

$$V(t, x; \pi) \equiv V(x; \pi) := \mathbb{E}^\mathbb{P} \left[ \int_0^\infty e^{-\beta s} [r(X_s^\pi, a_s^\pi) + \gamma p(X_s^\pi, a_s^\pi, \pi(\cdot \mid X_s^\pi))] ds \mid X_0^\pi = x \right]; \quad (37)$$

513 which is the state-value function under the policy  $\pi$ ,  $V(x; \pi)$ , in (4), and which, in turn, leads to the  
514 performance function  $\eta(\pi)$  in (6). Moreover, recall the main task of the continuous RL is to find (or  
515 approximate)  $\eta^* = \max_{\pi} \eta(\pi)$ , where max is over all admissible policies.

## 516 A.2 Controlled SDE and the HJ Equation

517 Note that the exploratory state dynamics in (34) is governed by a general Itô process. It is sometimes  
518 more convenient to consider an equivalent SDE representation—in the sense that its (weak) solution  
519 has the same distribution as the Itô process in (34) at each fixed time  $t$ . It is known ([58]) that when  
520  $n = m = 1$ , the marginal distribution of  $\{X_s^\pi, s \geq 0\}$  agrees with that of the solution to the SDE,  
521 denoted by  $\{\tilde{X}_s, s \geq 0\}$ :

$$d\tilde{X}_s = \tilde{b}(\tilde{X}_s, \pi(\cdot \mid \tilde{X}_s)) ds + \tilde{\sigma}(\tilde{X}_s, \pi(\cdot \mid \tilde{X}_s)) d\tilde{B}_s, \quad \tilde{X}_0 = x,$$

522 where  $\tilde{b}(x, \pi(\cdot)) = \int_{\mathcal{A}} b(x, a)\pi(a)da$  and  $\tilde{\sigma}(x, \pi(\cdot)) = \sqrt{\int_{\mathcal{A}} \sigma^2(x, a)\pi(a)da}$ . This result is easily  
523 extended to arbitrary  $n, m$ , thanks to [7, Corollary 3.7], with the precise statement presented below  
524 (assuming  $n = m$  for ease of exposition).

525 **Theorem 6.** *Assume that for a policy  $\pi$  and for every  $x$ ,*

$$\int_{\mathcal{A}} \sigma^2(x, a)\pi(a)da \in \mathbb{R}^{n \times n},$$

526 *is positive definite. Then there exists a filtered probability space  $(\tilde{\Omega}, \tilde{\mathcal{F}}, \{\tilde{\mathcal{F}}_t\}_{t \geq 0}, \tilde{\mathbb{P}})$  that supports  
527 a continuous  $\mathbb{R}^n$ -valued adapted process  $\tilde{X}$  and an  $n$ -dimensional Brownian motion  $\tilde{B}$  satisfying*

$$d\tilde{X}_s = \tilde{b}(\tilde{X}_s, \pi(\cdot \mid \tilde{X}_s)) ds + \tilde{\sigma}(\tilde{X}_s, \pi(\cdot \mid \tilde{X}_s)) d\tilde{B}_s, \quad \tilde{X}_0 = x, \quad (38)$$

528 *where*

$$\tilde{b}(x, \pi(\cdot)) = \int_{\mathcal{A}} b(x, a)\pi(a)da, \quad \tilde{\sigma}(x, \pi(\cdot)) = \left( \int_{\mathcal{A}} \sigma^2(x, a)\pi(a)da \right)^{\frac{1}{2}}.$$

529 *For each  $s \geq 0$ , the distribution of  $\tilde{X}_s$  under  $\tilde{\mathbb{P}}$  agrees with that of  $X_s^\pi$  under  $\mathbb{P}$  defined in (34).*

530 As a consequence, the state value function in (37) is identical to

$$V(x; \pi) = \mathbb{E} \left[ \int_0^\infty e^{-\beta s} \int_{\mathcal{A}} [r(\tilde{X}_s, a) + \gamma p(\tilde{X}_s, a, \pi(\cdot \mid \tilde{X}_s))] \pi(a \mid \tilde{X}_s) da ds \mid \tilde{X}_0 = x \right].$$

531 Also define

$$\tilde{r}(x, \pi) = \int_{\mathcal{A}} r(x, a) \pi(a|s) da, \quad \tilde{p}(x, \pi) = \int_{\mathcal{A}} p(x, a, \pi) \pi(a|x) da,$$

532 so we can simplify the value function to

$$V(x; \pi) = \mathbb{E} \left[ \int_0^\infty e^{-\beta s} \left[ \tilde{r}(\tilde{X}_s, \pi) + \gamma \tilde{p}(\tilde{X}_s^\pi, \pi(\cdot | \tilde{X}_s)) \right] ds \mid \tilde{X}_0 = x \right]. \quad (39)$$

533 Following the principle of optimality,  $V$  then satisfies the HJ equation:

$$\beta V(x; \pi) - \tilde{b}(x, \pi) \cdot \nabla V(x; \pi) - \frac{1}{2} \tilde{\sigma}^2(x, \pi) \circ \nabla^2 V(x; \pi) - \tilde{r}(x, \pi) - \gamma \tilde{p}(x, \pi) = 0. \quad (40)$$

534 To guarantee that the HJ equation in (40) characterizes the state-value function in (39), we need

535 **Assumption 3.** *Assume the following conditions hold:*

536 (i)  $b, \sigma, r, p$  are all continuous functions in their respective arguments.

537 (ii)  $\tilde{b}, \tilde{r}, \tilde{p}$  are uniformly Lipschitz continuous in  $x$ , i.e., there exists a constant  $C > 0$  such that for

538  $\varphi \in \{b, r\}$ ,

$$\|\varphi(x, a) - \varphi(x', a)\|_2 \leq C \|x - x'\|_2, \quad \text{for all } a \in \mathcal{A}, x, x' \in \mathbb{R}^n,$$

539 and

$$|p(x, a, \pi) - p(x', a, \pi)| \leq C \|x - x'\|_2, \quad \text{for all } a \in \mathcal{A}, \pi \in \mathcal{P}(\mathcal{A}), x, x' \in \mathbb{R}^n.$$

540 (iii)  $\tilde{\sigma}$  is globally bounded, i.e., there exist  $0 < \sigma_0 < \bar{\sigma}_0$  such that

$$\sigma_0^2 \cdot I \leq \tilde{\sigma}^2(x, a) \leq \bar{\sigma}_0^2 \cdot I, \quad \text{for all } a \in \mathcal{A}, x \in \mathbb{R}^n.$$

541 (iv) the SDE (38) has a weak solution which is unique in distribution.

542 (v)  $\pi(a|x)$  is measurable in  $(x, a)$  and is uniformly Lipschitz continuous in  $x$ , i.e., there exists a

543 constant  $C > 0$  such that

$$\int_{\mathcal{A}} |\pi(a|x) - \pi(a|x')| da \leq C \|x - x'\|_2, \quad \text{for all } x, x' \in \mathbb{R}^n.$$

544 **Theorem 7.** *Under Assumption 3, the state-value function in (39) is the unique (subquadratic)*

545 *viscosity solution to the HJ equation in (40).*

546 *Proof.* By [56, Section 3.1], the HJ equation in (40) has a unique (subquadratic) viscosity solution

547 under the conditions (i)-(iii). Further by [21, Lemma 2], the viscosity solution is the state-value

548 function.  $\square$

## 549 Appendix B Proofs of Main Results (in §3)

### 550 B.1 Proof of Theorem 2

551 Recall in the proof sketch of the Theorem in §3, we have defined the operator  $\mathcal{L}^\pi : C^2(\mathbb{R}^n) \mapsto C(\mathbb{R}^n)$

552 as

$$(\mathcal{L}^\pi \varphi)(x) := -\beta \varphi(x) + \tilde{b}(x, \pi) \cdot \nabla \varphi(x) + \frac{1}{2} \tilde{\sigma}^2(x, \pi) \circ \nabla^2 \varphi(x),$$

553 which leads to the following characterization of the HJ equation:

$$-\mathcal{L}^\pi V(x; \pi) = \tilde{r}(x, \pi) + \gamma \tilde{p}(x, \pi). \quad (41)$$

554 We need the following two lemmas concerning the operator  $\mathcal{L}^\pi$ .

555 **Lemma 8.** *For any  $\varphi \in C^2(\mathbb{R}^n)$ , we have*

$$\int_{\mathbb{R}^n} d_x^\pi(y) (-\mathcal{L}^\pi \varphi)(y) dy = \varphi(x).$$

556 *Proof.* The left hand side of the above equation is

$$\begin{aligned}
&= \mathbb{E} \int_0^\infty e^{-\beta s} \left( \beta \varphi(\tilde{X}_s^\pi) - \tilde{b}(\tilde{X}_s^\pi, \pi) \frac{\partial \varphi}{\partial x}(\tilde{X}_s^\pi) - \frac{1}{2} \tilde{\sigma}(\tilde{X}_s^\pi, \pi)^2 \frac{\partial^2 \varphi}{\partial x^2}(\tilde{X}_s^\pi) \right) ds \\
&= \mathbb{E} \int_0^\infty e^{-\beta s} \left[ \left( \beta \varphi(\tilde{X}_s^\pi) - \tilde{b}(\tilde{X}_s^\pi, \pi) \frac{\partial \varphi}{\partial x}(\tilde{X}_s^\pi) - \frac{1}{2} \tilde{\sigma}(\tilde{X}_s^\pi, \pi)^2 \frac{\partial^2 \varphi}{\partial x^2}(\tilde{X}_s^\pi) \right) ds - \tilde{\sigma}(\tilde{X}_s^\pi, \pi) \frac{\partial \varphi}{\partial x}(\tilde{X}_s^\pi) dB_s \right] \\
&= \mathbb{E} \int_0^\infty d \left( -e^{-\beta s} \varphi(\tilde{X}_s^\pi) \right) \\
&= \lim_{s \rightarrow \infty} \left( -e^{-\beta s} \varphi(\tilde{X}_s^\pi) \right) + \varphi(\tilde{X}_0^\pi) \\
&= \varphi(x),
\end{aligned}$$

557 where the first equality follows from the definition of the occupation time and the third equality from  
558 Itô's formula.  $\square$

559 **Lemma 9.** Let  $\pi, \hat{\pi}$  be two feedback policies. We have

$$(\mathcal{L}^{\hat{\pi}} - \mathcal{L}^\pi)V(x; \pi) + \tilde{r}(x, \hat{\pi}) - \tilde{r}(x, \pi) - \gamma \tilde{p}(x, \pi) = \int_{\mathcal{A}(x)} \hat{\pi}(a | x) q(x, a; \pi) da. \quad (42)$$

560 *Proof.* By definition of  $q(x, a; \pi)$  in (11), we have

$$\begin{aligned}
\text{RHS} &= \int_{\mathcal{A}(x)} \hat{\pi}(a | x) \left( \mathcal{H}^a \left( x, \frac{\partial V}{\partial x}(x; \pi), \frac{\partial^2 V}{\partial x^2}(x; \pi) \right) - \beta V(x; \pi) \right) da \\
&= \int_{\mathcal{A}(x)} \hat{\pi}(a | x) \left( b(x, a) \cdot \frac{\partial V}{\partial x}(x; \pi) + \frac{1}{2} \sigma^2(x, a) \circ \frac{\partial^2 V}{\partial x^2}(x; \pi) + r(x, a) - \beta V(x; \pi) \right) da \\
&= \tilde{r}(x, \hat{\pi}) + \mathcal{L}^{\hat{\pi}}V^\pi(x) \\
&= \tilde{r}(x, \hat{\pi}) - \tilde{r}(x, \pi) - \gamma \tilde{p}(x, \pi) + \mathcal{L}^{\hat{\pi}}V^\pi(x) - \mathcal{L}^\pi V^\pi(x) \\
&= \text{LHS}.
\end{aligned}$$

561  $\square$

*Proof of Theorem 2.* Note that in (13), the equation to be proven, the right hand side can be written as  $\int_{\mathbb{R}} d_\mu^{\hat{\pi}}(y) f(x; \pi, \hat{\pi}) dy$ , with

$$f(x; \pi, \hat{\pi}) := \int_{\mathcal{A}} \hat{\pi}(a | x) (q(x, a; \pi) + \gamma p(x, a, \hat{\pi})) da.$$

562 From Lemma 9, we have

$$f(x; \pi, \hat{\pi}) = (\mathcal{L}^{\hat{\pi}} - \mathcal{L}^\pi)V(x; \pi) + \tilde{r}(x, \hat{\pi}) + \gamma \tilde{p}(x, \hat{\pi}) - \tilde{r}(x, \pi) - \gamma \tilde{p}(x, \pi). \quad (43)$$

563 On the other hand, for the left hand side of (13), we have

$$\eta(\pi) = \int_{\mathbb{R}^n} V(y; \pi) \mu(dy) = \int_{\mathbb{R}^n} d_\mu^{\hat{\pi}}(y) (-\mathcal{L}^{\hat{\pi}})V(y; \pi) dy, \quad (44)$$

564 with the second equality following from Lemma 8; and

$$\eta(\hat{\pi}) = \int_{\mathbb{R}} d_\mu^{\hat{\pi}}(y) [\tilde{r}(y, \hat{\pi}) + \gamma \tilde{p}(y, \hat{\pi})] dy, \quad (45)$$

565 following the definition of the discounted expected occupation time; moreover, from (41), we have

$$0 = \int_{\mathbb{R}} d_\mu^{\hat{\pi}}(y) [(-\mathcal{L}^\pi)V(y; \pi) - \tilde{r}(y, \pi) - \gamma \tilde{p}(y, \pi)] dy. \quad (46)$$

566 Hence, combining the last three equations (44,45,46), we have

$$\eta(\hat{\pi}) - \eta(\pi) = \int_{\mathbb{R}} d_\mu^{\hat{\pi}}(y) [(\mathcal{L}^{\hat{\pi}} - \mathcal{L}^\pi)V(y; \pi) + \tilde{r}(y, \hat{\pi}) + \gamma \tilde{p}(y, \hat{\pi}) - \tilde{r}(y, \pi) - \gamma \tilde{p}(y, \pi)] dy. \quad (47)$$

567 Thus, we have shown LHS=RHS in (13).  $\square$

568 **B.2 Proof of Theorem 3**

569 *Proof.* It suffices to show the integral version of the theorem:

$$\nabla_{\theta} (\eta(\pi^{\theta})) |_{\theta=0} = \int_{\mathbb{R}^n} d_{\mu}^{\pi^{\theta}}(x) \left[ \int_{\mathcal{A}} \nabla_{\theta} \pi^{\theta}(a | x) (q(x, a; \pi^{\theta}) + \gamma p(x, a, \pi^{\theta})) + \gamma \cdot \pi^{\theta}(a | x) \nabla_{\theta} p(x, a, \pi^{\theta}) da \right] dx. \quad (48)$$

570 As before, we simplify notation by denoting  $\eta(\pi^{\theta})$  as  $\eta(\theta)$  and  $d^{\pi^{\theta}}$  as  $d^{\theta}$ . Then, by Theorem 2), we  
571 have

$$\eta(\theta + \delta\theta) - \eta(\theta) = \int_{\mathbb{R}^n} d_{\mu}^{\theta+\delta\theta}(x) \left[ \int_{\mathcal{A}} \pi^{\theta+\delta\theta}(a | x) (q(x, a; \theta) + \gamma p(x, a, \theta + \delta\theta)) da \right] dx. \quad (49)$$

572 Denote

$$f(\delta\theta) = \int_{\mathcal{A}} \pi^{\theta+\delta\theta}(a | x) (q(x, a; \theta) + \gamma p(x, a, \theta + \delta\theta)) da.$$

573 Note that  $f(0) = 0$ , which follows from

$$\begin{aligned} f(0) &= \int_{\mathcal{A}} \pi^{\theta}(a | x) (q(x, a; \theta) + \gamma p(x, a, \theta)) da \\ &= \int_{\mathcal{A}} \pi^{\theta}(a | x) \left( \mathcal{H}^a(x, \frac{\partial V}{\partial x}(x; \pi), \frac{\partial^2 V}{\partial x^2}(x; \pi)) - \beta V(x; \pi) + \gamma p(x, a, \theta) \right) da \\ &= -\beta V(x; \pi) + \tilde{b}(x, \pi) \cdot \nabla V(x; \pi) + \frac{1}{2} \tilde{\sigma}^2(x, \pi) \circ \nabla^2 V(x; \pi) + \tilde{r}(x, \pi) + \gamma \tilde{p}(x, \pi) \\ &= 0. \end{aligned}$$

574 Thus,

$$\begin{aligned} \eta(\theta + \delta\theta) - \eta(\theta) &= \langle d_{\mu}^{\theta+\delta\theta}, f(\delta\theta) \rangle \\ &= \langle d_{\mu}^{\theta+\delta\theta}, f(\delta\theta) \rangle - \langle d_{\mu}^{\theta+\delta\theta}, f(0) \rangle \\ &= \langle d_{\mu}^{\theta+\delta\theta}, f(\delta\theta) - f(0) \rangle \\ &= \langle d_{\mu}^{\theta+\delta\theta} - d_{\mu}^{\theta}, f(\delta\theta) - f(0) \rangle + \langle d_{\mu}^{\theta}, f(\delta\theta) - f(0) \rangle. \end{aligned}$$

575 Dividing both sides by  $\delta\theta$  completes the proof, as the first term on the last line above is of higher  
576 order than  $\delta\theta$ .  $\square$

577 **B.3 Proofs of Lemma 4 and Theorem 5**

578 We need a lemma for the perturbation bounds.

579 **Lemma 10.** Assume that both  $\tilde{\sigma}^2(x, \hat{\pi}(\cdot))$  and  $\tilde{\sigma}^2(x, \pi(\cdot))$  are positive definite and

$$\tilde{\sigma}^2(x, \pi(\cdot)), \tilde{\sigma}^2(x, \hat{\pi}(\cdot)) \geq \sigma_0^2 \cdot I.$$

580 where  $\sigma_0 > 0$ , then we have that the difference between the square root matrix is bounded by

$$\|\tilde{\sigma}(x, \hat{\pi}) - \tilde{\sigma}(x, \pi)\|_2 \leq \frac{1}{2\sigma_0} \|\tilde{\sigma}^2(x, \hat{\pi}) - \tilde{\sigma}^2(x, \pi)\|_2.$$

581 If we also assume that the upper bounds, i.e.

$$\tilde{\sigma}^2(x, \pi(\cdot)), \tilde{\sigma}^2(x, \hat{\pi}(\cdot)) \leq \bar{\sigma}_0^2 \cdot I.$$

582 by some  $\bar{\sigma}_0 > \sigma_0 > 0$ , then we have

$$\|\tilde{\sigma}(x, \hat{\pi}) - \tilde{\sigma}(x, \pi)\|_2 \leq \frac{\bar{\sigma}_0}{2\sigma_0} \|\hat{\pi} - \pi\|_1^{\frac{1}{2}}.$$

583 *Proof.* Consider a normalized vector  $x$  with  $\|x\|_2 = 1$  is an eigenvector of  $A^{\frac{1}{2}} - B^{\frac{1}{2}}$  with eigenvalue  
584  $\mu$  then

$$\begin{aligned} x^T (A - B)x &= x^T (A^{\frac{1}{2}} - B^{\frac{1}{2}}) A^{\frac{1}{2}} x + x^T B^{\frac{1}{2}} (A^{\frac{1}{2}} - B^{\frac{1}{2}}) x \\ &= \mu x^T (A^{\frac{1}{2}} + B^{\frac{1}{2}}) x. \end{aligned}$$

585 thus, if  $A, B \geq \sigma_0^2 I$ , this implies

$$\mu \leq \frac{|x^T(A-B)x|}{x^T(A^{\frac{1}{2}} + B^{\frac{1}{2}})x} \leq \|A-B\|_2 \cdot \lambda_{\min}(A^{\frac{1}{2}} + B^{\frac{1}{2}})^{-1} \leq \|A-B\|_2 / (2\sigma_0).$$

586 Furthermore, note that

$$\tilde{\sigma}^2(x, \hat{\pi}) - \tilde{\sigma}^2(x, \pi) = \int_{\mathcal{A}} \sigma^2(x, a)(\tilde{\pi}(a|x) - \pi(a|x)) da.$$

587 so

$$\|\tilde{\sigma}^2(x, \hat{\pi}) - \tilde{\sigma}^2(x, \pi)\|_2 \leq \bar{\sigma}_0^2 \int_{\mathcal{A}} |\tilde{\pi}(a|x) - \pi(a|x)| da = \bar{\sigma}_0^2 \cdot \|\tilde{\pi}(a|x) - \pi(a|x)\|_1.$$

588

□

589 *Proof* (of Lemma 4). Consider the Wasserstein-2 distance  $W_2(\mu, \nu)$  between distribution  $\mu$  and  $\nu$  as

$$W_2(\mu, \nu) = \left( \inf_{\gamma \in \Gamma(\mu, \nu)} \mathbf{E}_{(x,y) \sim \gamma} \|x - y\|_2^2 \right)^{1/2},$$

590 where  $\Gamma(\mu, \nu)$  is the set all probability measures on the product space  $\mathbb{R}^n \times \mathbb{R}^n$  with the marginal  
591 distributions being  $\mu$  and  $\nu$ , and  $\|\cdot\|_2$  is the standard Euclidean distance. Denote

$$\bar{d}_\mu^\pi := \beta d_\mu^\pi.$$

592 We want to get an upper bound on  $W_2(\bar{d}_\mu^\pi, \bar{d}_\mu^{\hat{\pi}})$  in terms of the distance between two policies  $\pi$  and  $\hat{\pi}$ .

593 Consider a specific coupling  $(X_t, Y_t)$  below:

$$\begin{cases} dX_s = \tilde{b}(X_s, \pi(\cdot | X_s)) ds + \tilde{\sigma}(X_s, \pi(\cdot | X_s)) dB_s, \\ dY_s = \tilde{b}(Y_s, \hat{\pi}(\cdot | Y_s)) ds + \tilde{\sigma}(Y_s, \hat{\pi}(\cdot | Y_s)) dB_s. \end{cases} \quad (50)$$

594 with  $X_0 = Y_0$ , which leads to a joint distribution over  $\mathbb{R}^n \times \mathbb{R}^n$ :

$$\tilde{\gamma} := \left\{ \tilde{p}(x, y) = \int_0^\infty \frac{1}{\beta} e^{-\beta t} f_{(X_t, Y_t)}(x, y) dt \right\}.$$

595 Hence,

$$W_2^2(\bar{d}_\mu^\pi, \bar{d}_\mu^{\hat{\pi}}) \leq \mathbb{E}_{(x,y) \sim \tilde{\gamma}} \|x - y\|_2^2 = \int_0^\infty \frac{1}{\beta} e^{-\beta s} \mathbb{E} \|X_s - Y_s\|_2^2 ds. \quad (51)$$

596 It then boils down to estimating  $\mathbb{E} \|X_s - Y_s\|_2^2$ . By Itô's formula,

$$\begin{aligned} d\|X_s - Y_s\|_2^2 &= 2(X_s - Y_s)^\top \left[ (\tilde{b}(X_s, \pi) - \tilde{b}(Y_s, \hat{\pi})) ds + (\tilde{\sigma}(X_s, \pi) - \tilde{\sigma}(Y_s, \hat{\pi})) dB_s \right] \\ &\quad + \text{Tr} [(\tilde{\sigma}(X_s, \pi) - \tilde{\sigma}(Y_s, \hat{\pi}))^2] ds. \end{aligned}$$

597 Taking expectation on both sides yields

$$\frac{d}{ds} \mathbb{E} \|X_s - Y_s\|_2^2 = 2 \underbrace{\mathbb{E} \left[ (X_s - Y_s)^\top (\tilde{b}(X_s, \pi) - \tilde{b}(Y_s, \hat{\pi})) ds \right]}_{(A)} + \underbrace{\text{Tr} \left[ \mathbb{E} (\tilde{\sigma}(X_s, \pi) - \tilde{\sigma}(Y_s, \hat{\pi}))^2 \right]}_{(B)}, \quad (52)$$

598 with

$$\begin{aligned} (A) &= \mathbb{E} \left[ (X_s - Y_s)^\top (\tilde{b}(X_s, \pi) - \tilde{b}(Y_s, \pi)) ds \right] + \mathbb{E} \left[ (X_s - Y_s)^\top (\tilde{b}(Y_s, \pi) - \tilde{b}(Y_s, \hat{\pi})) ds \right] \\ &\leq C_{\tilde{b}} \cdot \mathbb{E} \|X_s - Y_s\|_2^2 + \frac{1}{2} \mathbb{E} \|X_s - Y_s\|_2^2 + \frac{1}{2} \mathbb{E} \|\tilde{b}(Y_s, \pi) - \tilde{b}(Y_s, \hat{\pi})\|_2^2 \\ &\leq (C_{\tilde{b}} + \frac{1}{2}) \cdot \mathbb{E} \|X_s - Y_s\|_2^2 + \frac{1}{2} \|\tilde{b}(\cdot, \pi) - \tilde{b}(\cdot, \hat{\pi})\|_{2, \infty}^2; \end{aligned}$$

599 and

$$\begin{aligned} (B) &= \mathbb{E} \|\tilde{\sigma}(X_s, \pi) - \tilde{\sigma}(Y_s, \hat{\pi})\|_F^2 \\ &\leq 2\mathbb{E} \|\tilde{\sigma}(X_s, \pi) - \tilde{\sigma}(Y_s, \pi)\|_F^2 + 2\mathbb{E} \|\tilde{\sigma}(Y_s, \pi) - \tilde{\sigma}(Y_s, \hat{\pi})\|_F^2 \\ &\leq 2C_{\tilde{\sigma}}^2 \cdot \mathbb{E} \|X_s - Y_s\|_2^2 + 2 \sup_x \|\tilde{\sigma}(x, \pi) - \tilde{\sigma}(x, \hat{\pi})\|_F^2 \\ &:= 2C_{\tilde{\sigma}}^2 \cdot \mathbb{E} \|X_s - Y_s\|_2^2 + 2\|\tilde{\sigma}(\cdot, \pi) - \tilde{\sigma}(\cdot, \hat{\pi})\|_{F, \infty}^2. \end{aligned}$$

600 Combining the above, we get

$$\frac{d}{ds} \mathbb{E} \|X_s - Y_s\|_2^2 \leq \underbrace{(2C_{\tilde{b}} + 1 + 2C_{\tilde{\sigma}}^2)}_{C_{\tilde{b}, \tilde{\sigma}}} \mathbb{E} \|X_s - Y_s\|_2^2 + \underbrace{\|\tilde{b}(\cdot, \pi) - \tilde{b}(\cdot, \hat{\pi})\|_{2, \infty}^2 + 2\|\tilde{\sigma}(\cdot, \pi) - \tilde{\sigma}(\cdot, \hat{\pi})\|_{F, \infty}^2}_{C(\pi, \hat{\pi})}.$$

601 By Grönwall's inequality, we have

$$\mathbb{E} \|X_t - Y_t\|_2^2 \leq \frac{C(\pi, \hat{\pi})}{C_{\tilde{b}, \tilde{\sigma}}} (e^{C_{\tilde{b}, \tilde{\sigma}} t} - 1). \quad (53)$$

602 Substituting back into (51), we obtain

$$W_2^2(\bar{d}_\mu^\pi, \bar{d}_\mu^{\hat{\pi}}) \leq \frac{C(\pi, \hat{\pi})}{C_{\tilde{b}, \tilde{\sigma}}} \int_0^\infty \frac{1}{\beta} e^{-\beta s} (e^{C_{\tilde{b}, \tilde{\sigma}} s} - 1) ds.$$

603 Thus, if  $\beta > C_{\tilde{b}, \tilde{\sigma}}$ , we have

$$W_2(\bar{d}_\mu^\pi, \bar{d}_\mu^{\hat{\pi}}) \leq \frac{C(\pi, \hat{\pi})}{C_{\tilde{b}, \tilde{\sigma}}(\beta - C_{\tilde{b}, \tilde{\sigma}})\beta}.$$

604 Concerning the term  $C(\pi, \hat{\pi})$ , we have

$$\|\tilde{b}(\cdot, \pi) - \tilde{b}(\cdot, \hat{\pi})\|_{2, \infty} = \sup_x \|\tilde{b}(x, \pi) - \tilde{b}(x, \hat{\pi})\|_2 \leq \sup_x \|\hat{\pi}(\cdot|x) - \pi(\cdot|x)\|_1 \cdot \sup_{x, a} |b(x, a)|,$$

605 and

$$\|\tilde{\sigma}(\cdot, \pi) - \tilde{\sigma}(\cdot, \hat{\pi})\|_{F, \infty} = \sup_x \|\tilde{\sigma}(x, \pi) - \tilde{\sigma}(x, \hat{\pi})\|_F \leq \sqrt{n} \frac{\bar{\sigma}_0}{2\sigma_0} \sup_x \|\hat{\pi}(\cdot|x) - \pi(\cdot|x)\|_1^{\frac{1}{2}}.$$

606 Thus we have:

$$\begin{aligned} C(\pi, \hat{\pi}) &= \|\tilde{b}(\cdot, \pi) - \tilde{b}(\cdot, \hat{\pi})\|_{2, \infty}^2 + 2\|\tilde{\sigma}(\cdot, \pi) - \tilde{\sigma}(\cdot, \hat{\pi})\|_{F, \infty}^2 \\ &\leq \left( \sup_{x, a} |b(x, a)|^2 + \frac{d \cdot \bar{\sigma}_0^2}{2\sigma_0^2} \right) \max \left( \sup_x \|\hat{\pi}(\cdot|x) - \pi(\cdot|x)\|_1, \sup_x \|\hat{\pi}(\cdot|x) - \pi(\cdot|x)\|_1^{\frac{1}{2}} \right) \end{aligned}$$

607 which proves our upper bound.  $\square$

608 *Proof* (of Theorem 5). We have that

$$\begin{aligned} |\eta^{\hat{\pi}} - L^\pi(\hat{\pi})| &= |\langle \bar{d}_\mu^{\hat{\pi}} - \bar{d}_\mu^\pi, f \rangle| = \frac{\|f\|_{\dot{H}^1}}{\beta} \left| \left\langle \bar{d}_\mu^{\hat{\pi}} - \bar{d}_\mu^\pi, \frac{f}{\|f\|_{\dot{H}^1}} \right\rangle \right| \\ &\leq \frac{K}{\beta} \|\bar{d}_\mu^{\hat{\pi}} - \bar{d}_\mu^\pi\|_{\dot{H}^{-1}} \leq \frac{K\sqrt{M}}{\beta} W_2(\bar{d}_\mu^{\hat{\pi}}, \bar{d}_\mu^\pi). \end{aligned} \quad (54)$$

609 where  $K := \sup_{\hat{\pi}} \|f\|_{\dot{H}^1} < \infty$  (more about  $K$  in the remarks below). Combining (54) with the  
610 estimate in (22) (of Lemma 4) yields the desired result in (23).  $\square$

611 *Remarks* (on  $K$ ). In the performance-difference bound developed above, we assume  $K$  is finite:

$$K := \|f\|_{\dot{H}^1} := \left( \int_{\mathbb{R}^n} |\nabla f(x)|^2 dx \right)^{\frac{1}{2}} < \infty,$$

612 where  $f(x; \pi, \hat{\pi}) := \int_{\mathcal{A}} \hat{\pi}(a|x) (q(x, a; \pi) + p(x, a, \hat{\pi})) da$ . The famous Poincaré inequality can  
613 provide a lower bound on this quantity; but we need an upper bound as well, i.e.,

$$K = \left( \int_{\mathbb{R}^n} |\nabla f(x)|^2 dx \right)^{\frac{1}{2}} \leq C \left( \int_{\mathbb{R}^n} |f(x)|^2 dx \right)^{\frac{1}{2}}.$$

614 This above is essentially a *reverse* Poincaré Inequality, which is not likely to hold (in particular, the  
615 existence of the constant  $C$ ).

<sup>1</sup>From this proof, it's evident that there's a  $\beta$  missing in the denominator on the RHS of (22). Consequently, the  $C(\mu, \pi, \hat{\pi})$  expression in Theorem 5 should have  $2\beta^2$  (instead of  $2\beta$ ) in the denominator. This correction will *not* affect the two numerical examples as both had set  $\beta = 1$  (as a hyper-parameter).

616 Should we indeed have a reverse Poincaré Inequality, then we can further bound  $f$  by

$$\begin{aligned}
 |f(x)| &= \left| \int_{\mathcal{A}} (\hat{\pi}(a | x) - \pi(a | x)) (q(x, a; \pi) + p(x, a, \hat{\pi})) da \right| \\
 &\leq \int_{\mathcal{A}} |\hat{\pi}(a | x) - \pi(a | x)| \cdot |q(x, a; \pi) + p(x, a, \hat{\pi})| da \\
 &\leq 2 \sup_a |q(x, a; \pi) + p(x, a, \hat{\pi})| D_{\text{TV}}(\pi(\cdot | x), \hat{\pi}(\cdot | x)),
 \end{aligned}$$

617 and

$$\begin{aligned}
 \left( \int_{\mathbb{R}^n} |f(x)|^2 dx \right)^{\frac{1}{2}} &\leq \left( \int_{\mathbb{R}^n} 4 \sup_a |q(x, a; \pi) + p(x, a, \hat{\pi})|^2 D_{\text{TV}}^2(\pi(\cdot | x), \hat{\pi}(\cdot | x)) dx \right)^{\frac{1}{2}} \\
 &\leq \left( \int_{\mathbb{R}^n} 2 \sup_a |q(x, a; \pi) + p(x, a, \hat{\pi})|^2 dx \right)^{\frac{1}{2}} \sqrt{\sup_x D_{\text{KL}}(\pi(\cdot | x), \hat{\pi}(\cdot | x))},
 \end{aligned}$$

618 where the second inequality is from Pinsker's inequality. This way, we would have recovered a  
 619 similar bound as in the discrete RL. Since we do not have the reverse Poincaré inequality, however,  
 620 we have to assume that  $K$  is finite.

## 621 Appendix C Algorithms

### 622 C.1 Performance of CPPO with Square-root KL and Linear KL

623 Here we present a detailed version of the CPPO algorithm. For two probability distributions  $P$  and  $Q$   
 624 over the action space with density functions  $p$  and  $q$  correspondingly, the KL-divergence between  
 625 these two is defined as:

$$D_{\text{KL}}(P\|Q) = \int_{\mathcal{A}} \log\left(\frac{q(a)}{p(a)}\right)q(a)da,$$

626 Denote  $D_{\text{KL}}(\theta, \theta_k) := \mathbb{E}_{x \sim d_{\mu}^{\theta_k}} D_{\text{KL}}(\pi_{\theta}(\cdot|x)\|\pi_{\theta_k}(\cdot|x))$ , to distinguish it from  $\bar{D}_{\text{KL}}(\theta\|\theta_k) :=$   
 627  $\mathbb{E}_{x \sim d_{\mu}^{\theta_k}} \sqrt{D_{\text{KL}}(\pi_{\theta}(\cdot|x)\|\pi_{\theta_k}(\cdot|x))}$  which was used in CPPO Algorithm in 2.

628 Note that bounding the performance difference by the linear KL-divergence  $D_{\text{KL}}(\theta, \theta_k)$ , instead of  
 629 its square-root counterpart  $\bar{D}_{\text{KL}}(\theta\|\theta_k)$ , will generally require stronger conditions (which may be  
 630 difficult to satisfy). For completeness, we present the following algorithm, the CPPO with linear  
 631 KL-divergence:

---

#### Algorithm 3 CPPO: PPO with adaptive penalty constant (linear KL-divergence)

---

**Input:** Policy parameters  $\theta_0$ , critic net parameters  $\phi_0$

- 1: **for**  $k = 0, 1, 2, \dots$  until  $\theta_k$  converge **do**
- 2:   Collect a truncated trajectory  $\{X_{t_i}, a_{t_i}, r_{t_i}, p_{t_i}\}, i = 1, \dots, N$  from the environment using  $\pi_{\theta_k}$ .
- 3:   **for**  $i = 0, \dots, N - 1$  **do:** Update the critic parameters as in (8)
- 4:   **for**  $j = 1, \dots, J$  **do:** Draw i.i.d.  $\tau_j$  from  $\exp(\beta)$ , round  $\tau_j$  to the largest multiple of  $\delta_t$  no larger than it, and compute the GAE estimator of  $q(X_{\tau_j}, a_{\tau_j})$

$$\tilde{q}(X_{\tau_j}, a_{\tau_j}) := (r_{\tau_j} \delta_t + e^{-\beta \delta_t} V(X_{\tau_j + \delta_t}) - V(X_{\tau_j})) / \delta_t.$$

- 5:   Compute policy update (by taking a fixed  $s$  steps of gradient descent)

$$\theta_{k+1} = \arg \max_{\theta} L^{\theta_k}(\theta) - C_{\text{penalty}}^k D_{\text{KL}}(\theta, \theta_k).$$

- 6:   **if**  $D_{\text{KL}}(\theta_{k+1}, \theta_k) \geq (1 + \epsilon)\delta$ , **then**  $C_{\text{penalty}}^{k+1} = 2C_{\text{penalty}}^k$ .
  - 7:   **else if**  $D_{\text{KL}}(\theta_{k+1}, \theta_k) \leq \delta/(1 + \epsilon)$ , **then**  $C_{\text{penalty}}^{k+1} = C_{\text{penalty}}^k/2$ .
- 

632 A comparison between the above and Algorithm 2 (using square-root KL divergence) is presented in  
 633 §D.3 below, which clearly illustrates the advantage of square-root KL divergence.

### 634 C.2 KL-divergence

635 We elaborate here on the KL-divergence between the current policy and the optimal policy, along  
 636 with the entropy regularizer. By the performance difference formula, we have

$$\eta(\pi) - \eta(\pi^*) = \int_{\mathbb{R}^n} d_{\mu}^{\pi}(x) \left[ \int_{\mathcal{A}} \pi(a|x) (q(x, a; \pi^*) - \gamma \log(\pi(a))) da \right] dx.$$

637 Notice that by the definition of KL-divergence we defined before, we have

$$D_{\text{KL}}(\pi^*(\cdot|x)\|\pi(\cdot|x)) = \int_{\mathcal{A}} \log\left(\frac{\pi(a|x)}{\pi^*(a|x)}\right)\pi(a|x)da.$$

638 Similar as the previous discussion of soft  $q$ -learning,  $\pi^*$  is optimal implies that

$$\pi^*(a|x) \propto \exp\left(\frac{q(x, a, \pi^*)}{\gamma}\right),$$

639 and the normalization constant is 1 can be proved through considering the exploratory HJB equation,  
 640 see [22, 56]. Thus

$$D_{\text{KL}}(\pi^*(\cdot|x)\|\pi(\cdot|x)) = \int_{\mathcal{A}} \log(\pi(a|x))\pi(a|x)da - \int_{\mathcal{A}} \frac{q(x, a, \pi^*)}{\gamma} \pi(a|x)da,$$

641 which leads to

$$\eta(\pi) - \eta(\pi^*) = -\gamma \cdot \mathbb{E}_{x \sim d_{\mu}^{\pi}} D_{\text{KL}}(\pi^*(\cdot|x) \parallel \pi(\cdot|x)).$$

642 This justifies our claim that the KL-divergence is essentially equivalent to the distance to the optimal  
643 performance.

644 **Appendix D Experiments**

645 **D.1 Example 1**

646 Recall, in the LQ control problem, the reward function is

$$r(x, a) = - \left( \frac{M}{2} x^2 + Rxa + \frac{N}{2} a^2 + Px + Qa \right),$$

647 with  $M \geq 0, N > 0, R, Q, P \in \mathbb{R}$  and  $R^2 < MN$ , and we adopt the entropy regularizer as

$$p(x, a, \pi) = -\log(\pi(a)).$$

648 Furthermore, suppose that the discount rate satisfies  $\beta > 2A + C^2 + \max\left(\frac{D^2 R^2 - 2NR(B+CD)}{N}, 0\right)$ .

649 The following results are readily derived from Theorem 4 of [58]. The value function of the optimal  
650 policy  $\pi^*$  is

$$V(x) = \frac{1}{2} k_2 x^2 + k_1 x + k_0, \quad x \in \mathbb{R},$$

651 where

$$k_2 := \frac{1}{2} \frac{(\rho - (2A + C^2))N + 2(B + CD)R - D^2 M}{(B + CD)^2 + (\rho - (2A + C^2))D^2}$$

$$- \frac{1}{2} \frac{\sqrt{((\rho - (2A + C^2))N + 2(B + CD)R - D^2 M)^2 - 4((B + CD)^2 + (\rho - (2A + C^2))D^2)(R^2 - MN)}}{(B + CD)^2 + (\rho - (2A + C^2))D^2},$$

$$k_1 := \frac{P(N - k_2 D^2) - QR}{k_2 B(B + CD) + (A - \rho)(N - k_2 D^2) - BR},$$

652 and

$$k_0 := \frac{(k_1 B - Q)^2}{2\rho(N - k_2 D^2)} + \frac{\gamma}{2\rho} \left( \ln \left( \frac{2\pi e \gamma}{N - k_2 D^2} \right) - 1 \right)$$

653 respectively. Moreover, the optimal feedback control is Gaussian, with density function

$$\pi^*(a; x) = \mathcal{N} \left( a \mid \frac{(k_2(B + CD) - R)x + k_1 B - Q}{N - k_2 D^2}, \frac{\gamma}{N - k_2 D^2} \right).$$

654 For a set of model parameters:  $A = -1, B = C = 0, D = 1, M = N = Q = 2, R = P = 1, \beta =$   
655  $1, \gamma = 0.1$ , following the formulas and the parameterized policy  $\pi_\theta(\cdot \mid x) = \mathcal{N}(\theta_1 x + \theta_2, \exp(\theta_3))$ ,  
656 and the corresponding value function  $V_\phi(x) = \frac{1}{2}\phi_2 x^2 + \phi_1 x + \phi_0$ , we can derive the optimal  
657 parameters:

$$\phi^* = [0.71914874, -0.10555128, -0.53518376],$$

658 and

$$\theta^* = [-0.39444872, -0.78889745, -1.40400944].$$

Table 1: Hyper-parameter values for Example 1

Alphabet	Description	Value
$T$	Trajectory Truncation Length	25
$\beta$	discount factor	1
$\delta_t$	time interval	0.005
$J$	batch size for sampling $\exp(\beta)$	100
$\alpha_1$	learning rate for policy iteration $k$	0.02 when $k \leq 50$ and $0.02 \times \log(\frac{50}{k})$ when $k > 50$
$\alpha_2$	learning rate for value iteration $k$	0.01 when $k \leq 50$ and $0.01 \times \log(\frac{50}{k})$ when $k > 50$
$K$	iteration threshold	2000
$s$	steps of gradient descent	10
$\delta$	radius	0.0002
$\epsilon$	tolerance level	0.5

659 **D.2 Example 2**

660 The model parameters are  $k = 0.01, \theta = 7, \eta = 0.1, \rho = 0.3, \sigma = 1, r_f = 0.01, \ell = 5$ . For both the  
 661 value function and the policy parameterization, we use a 3-layer neural network, and with the initial  
 662 parameters sampled from the uniform distribution over  $[-0.5, 0.5]$ . We use the tanh activation function  
 663 for the hidden layer.

Table 2: Hyperparameter values for Example 2

Alphabet	Description	Value
$T$	Trajectory Truncation Length	25
$\beta$	discount factor	1
$\delta_t$	time interval	0.005
$J$	batch size for sampling $\exp(\beta)$	100
$\alpha_1$	learning rate for policy iteration $k$	0.005 when $k \leq 50$ and $0.005 \times \log(\frac{50}{k})$ when $k > 50$
$\alpha_2$	learning rate for value iteration $k$	0.01 when $k \leq 50$ and $0.01 \times \log(\frac{50}{k})$ when $k > 50$
$K$	iteration threshold	200
$s$	steps of gradient descent	10
$\delta$	radius	0.025
$\epsilon$	tolerance level	0.5

664 **D.3 Performance of CPPO with Square-root KL and Linear KL**

665 We compare the performance of CPPO with square-root KL-divergence (denote as CPPO), and linear  
 666 KL-divergence (denoted as CPPO (nst) — non square-root) applied to the experiments in Example  
 1 and Example 2. Figure 4 compares the distance between the current policy parameters and the

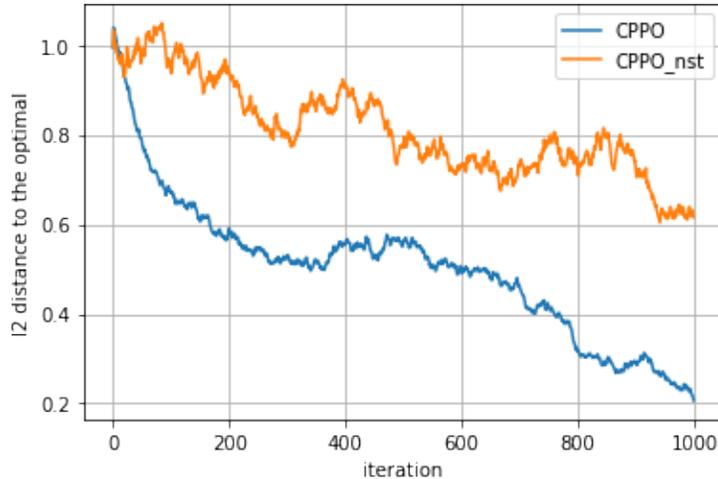


Figure 4: Performance of CPPO and CPPO (nst) to the Example 1

667 optimal parameters, with  $x$ -axis denoting the iteration times and  $y$ -axis denoting the  $L_2$  distance.  
 668 Figure 5 compares the current expected return, with  $x$ -axis denoting the iteration times and  $y$ -axis  
 669 denoting the current performance by taking the average of 100 times of Monte Carlo evaluation. In  
 670 both figures, the blue curve represents the algorithm with square-root KL-divergence as opposed to  
 671 the orange one corresponding to the linear version. Both figures clearly demonstrate the advantage  
 672 of the former. In particular, the linear version can suffer from getting stuck at the local optimum as  
 673 demonstrated in Example 1.  
 674

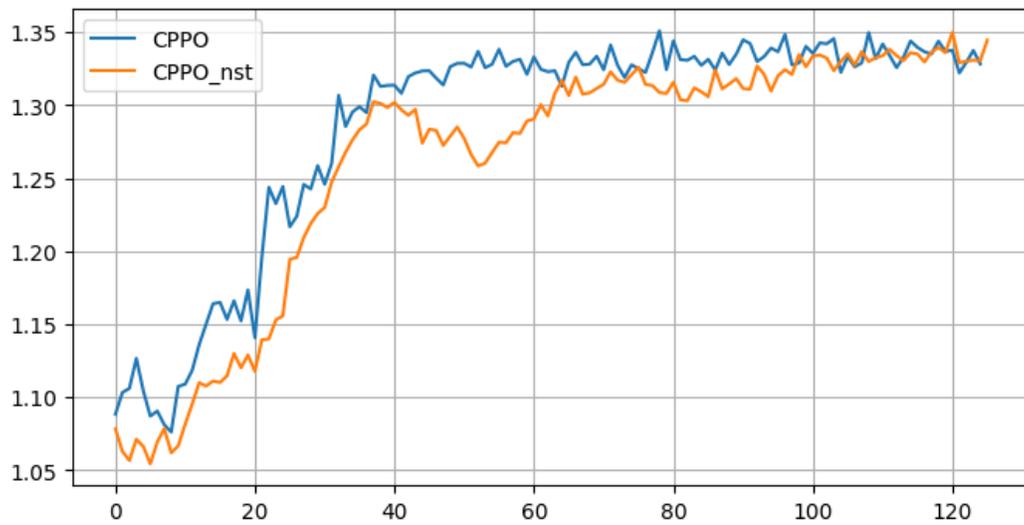


Figure 5: Performance of CPPO and CPPO (nst) to the Example 2