

1 A Mathematical Proofs

2 Here we provide proofs for Theorem 1 in the main paper.

3 **Theorem.** *The gap between the expected return of the model and the environment is bounded as:*

$$|J(\pi) - J^{\hat{P}}(\pi)| \leq \frac{2R_{max}}{(1-\gamma)^2} \left((2-\gamma)\epsilon_\pi + (1-\gamma) \sum_{t=1}^{\infty} \gamma^t \epsilon_t^m \right), \quad (1)$$

4 where $\epsilon_\pi := \max_s D_{TV}(\pi_D(\cdot|s) \|\pi(\cdot|s))$ denotes the policy distribution shift, $\epsilon_t^m :=$
 5 $\mathbb{E}_{s \sim \hat{P}_{t-1}(s,a;\pi)} [D_{TV}(\hat{P}_t(\cdot|s,a) \|\hat{P}_t(\cdot|s,a))]$ denotes the upper bound of one-step model prediction
 6 error at timestep t of the model rollout trajectory, $D_{TV}(p\|q)$ refers to the total variation between dis-
 7 tribution p and q , $R_{max} := \max_{s,a} R(s,a)$, and $\hat{P}_{t-1}(s,a;\pi)$ denotes the state-action distribution
 8 at t under \hat{P} and π .

Proof.

$$\begin{aligned} |J^{\hat{P}}(\pi) - J(\pi_D)| &= \left| \sum_{t=0}^{\infty} \gamma^t \sum_{s,a} \left(\hat{P}_t^\pi(s,a) - P_t^{\pi_D}(s,a) \right) R(s,a) \right| \\ &\leq R_{max} \sum_{t=0}^{\infty} \gamma^t \sum_{s,a} \left| \hat{P}_t^\pi(s,a) - P_t^{\pi_D}(s,a) \right| \\ &\leq 2R_{max} \sum_{t=0}^{\infty} \gamma^t D_{TV} \left(\hat{P}_t^\pi(s,a) \|\hat{P}_t^{\pi_D}(s,a) \right). \end{aligned} \quad (2)$$

9 Applying Lemma 1, we have:

$$D_{TV} \left(\hat{P}_t^\pi(s,a) \|\hat{P}_t^{\pi_D}(s,a) \right) \leq D_{TV} \left(\hat{P}_t^\pi(s) \|\hat{P}_t^{\pi_D}(s) \right) + \epsilon_\pi. \quad (3)$$

10 Similar to the proof of Lemma B.2 in [Janner et al., 2019], we have:

$$\begin{aligned} \left| \hat{P}_t^\pi(s) - P_t^{\pi_D}(s) \right| &= \left| \sum_{s'} \hat{P}^\pi(s_t = s | s') \hat{P}_{t-1}^\pi(s') - P^{\pi_D}(s_t = s | s') P_{t-1}^{\pi_D}(s') \right| \\ &\leq \sum_{s'} \left| \hat{P}^\pi(s_t = s | s') \hat{P}_{t-1}^\pi(s') - P^{\pi_D}(s_t = s | s') P_{t-1}^{\pi_D}(s') \right| \\ &= \sum_{s'} \left| \hat{P}^\pi(s | s') \hat{P}_{t-1}^\pi(s') - P^{\pi_D}(s | s') \hat{P}_{t-1}^\pi(s') + P^{\pi_D}(s | s') \hat{P}_{t-1}^\pi(s') - P^{\pi_D}(s | s') P_{t-1}^{\pi_D}(s') \right| \\ &\leq \sum_{s'} \hat{P}_{t-1}^\pi(s') \left| \hat{P}^\pi(s | s') - P^{\pi_D}(s | s') \right| + P^{\pi_D}(s | s') \left| \hat{P}_{t-1}^\pi(s') - P_{t-1}^{\pi_D}(s') \right| \\ &= \mathbb{E}_{s' \sim \hat{P}_{t-1}^\pi} \left[\left| \hat{P}^\pi(s | s') - P^{\pi_D}(s | s') \right| \right] + \sum_{s'} P^{\pi_D}(s | s') \left| \hat{P}_{t-1}^\pi(s') - P_{t-1}^{\pi_D}(s') \right|. \end{aligned} \quad (4)$$

11 Then, the term $D_{TV} \left(\hat{P}_t^\pi(s) \|\hat{P}_t^{\pi_D}(s) \right)$ can be bounded by:

$$\begin{aligned}
D_{TV} \left(\hat{P}_t^\pi(s) \| P_t^{\pi_D}(s) \right) &= \frac{1}{2} \sum_s \left| \hat{P}_t^\pi(s) - P_t^{\pi_D}(s) \right| \\
&= \frac{1}{2} \mathbb{E}_{s' \sim \hat{P}_{t-1}^\pi} \left[\sum_s \left| \hat{P}^\pi(s | s') - P^{\pi_D}(s | s') \right| \right] + \frac{1}{2} D_{TV} \left(\hat{P}_{t-1}^\pi(s') \| P_{t-1}^{\pi_D}(s') \right) \\
&= \frac{1}{2} \sum_{t'=1}^t \mathbb{E}_{s' \sim \hat{P}_{t'-1}^\pi} \left[\sum_s \left| \hat{P}^\pi(s | s') - P^{\pi_D}(s | s') \right| \right] \\
&= \frac{1}{2} \sum_{t'=1}^t \mathbb{E}_{s' \sim \hat{P}_{t'-1}^\pi} \left[\sum_s \left| \sum_a \hat{P}^\pi(s, a | s') - P^{\pi_D}(s, a | s') \right| \right] \\
&\leq \frac{1}{2} \sum_{t'=1}^t \mathbb{E}_{s' \sim \hat{P}_{t'-1}^\pi} \left[\sum_{s,a} \left| \hat{P}^\pi(s, a | s') - P^{\pi_D}(s, a | s') \right| \right] \\
&= \sum_{t'=1}^t \mathbb{E}_{s' \sim \hat{P}_{t'-1}^\pi} D_{TV} \left(\hat{P}^\pi(s, a | s') \| P^{\pi_D}(s, a | s') \right). \tag{5}
\end{aligned}$$

12 Again applying Lemma 1, we have:

$$D_{TV} \left(\hat{P}^\pi(s, a | s') \| P^{\pi_D}(s, a | s') \right) \leq \epsilon_\pi + \mathbb{E}_{a \sim \pi} D_{TV} \left(\hat{P}(s | s', a) \| P(s | s', a) \right). \tag{6}$$

13 Plugging Eq. (6) in Eq. (5) we have:

$$D_{TV} \left(\hat{P}_t^\pi(s) \| P_t^{\pi_D}(s) \right) \leq \sum_{t'=1}^t \epsilon_\pi + \epsilon_{t'}^m = t\epsilon_\pi + \sum_{t'=1}^t \epsilon_{t'}^m. \tag{7}$$

14 Plugging Eq. (7) in Eq. (3) we have:

$$D_{TV} \left(\hat{P}_t^\pi(s, a) \| P_t^{\pi_D}(s, a) \right) \leq t\epsilon_\pi + \sum_{t'=1}^t \epsilon_{t'}^m + \epsilon_\pi = (t+1)\epsilon_\pi + \sum_{t'=1}^t \epsilon_{t'}^m. \tag{8}$$

15 Plugging Eq. (8) in Eq. (2) we have:

$$\begin{aligned}
\left| J^{\hat{P}}(\pi) - J(\pi_D) \right| &\leq 2R_{max} \sum_{t=0}^{\infty} \gamma^t \left((t+1)\epsilon_\pi + \sum_{t'=1}^t \epsilon_{t'}^m \right) \\
&= 2R_{max} \left(\frac{\epsilon_\pi}{(1-\gamma)^2} + \frac{1}{(1-\gamma)} \sum_{t=1}^{\infty} \gamma^t \epsilon_t^m \right) \\
&= \frac{2R_{max}}{(1-\gamma)^2} \left(\epsilon_\pi + (1-\gamma) \sum_{t=1}^{\infty} \gamma^t \epsilon_t^m \right). \tag{9}
\end{aligned}$$

16 Therefore, the result in Eq. (1) can be derived:

$$\begin{aligned}
\left| J(\pi) - J^{\hat{P}}(\pi) \right| &\leq \left| J(\pi) - J(\pi_D) \right| + \left| J(\pi_D) - J^{\hat{P}}(\pi) \right| \\
&\leq \frac{2R_{max}\epsilon_\pi}{1-\gamma} + \frac{2R_{max}}{(1-\gamma)^2} \left(\epsilon_\pi + (1-\gamma) \sum_{t=1}^{\infty} \gamma^t \epsilon_t^m \right) \\
&= \frac{2R_{max}}{(1-\gamma)^2} \left((2-\gamma)\epsilon_\pi + (1-\gamma) \sum_{t=1}^{\infty} \gamma^t \epsilon_t^m \right).
\end{aligned}$$

17

□

18 **Lemma 1.** (TVD of joint distribution) Suppose that we have two distributions $P_1(x, y) =$
 19 $P_1(x)P_1(y|x)$ and $P_2(x, y) = P_2(x)P_2(y|x)$. We can bound the total variance difference of the joint
 20 as:

$$D_{TV}(P_1(x, y)||P_2(x, y)) \leq D_{TV}(P_1(x)||P_2(x)) + \mathbb{E}_{x \sim P_1}[D_{TV}(P_1(y|x)||P_2(y|x))].$$

21 Lemma 1 is proved in the MBPO paper, so we only provide the result here.

22 B Experimental Details

23 We describe some implementation details and hyperparameter settings below.

24 B.1 Implementation and Hyperparameter Settings

25 Our implementation is overall based on MBPO. The algorithm for policy learning, the actor-critic
 26 network architecture, the bootstrapped model ensemble technique and other details are all the same
 27 with those in MBPO. The only modified part is the model learning process which is the focus of
 28 the main paper. In P2P-MPC, the candidate action sequences are generated in a parallel manner to
 29 accelerate this process, at the cost of some extra memory cost. The number of candidate sequences
 30 is set to 4 for Hopper, Ant and Humanoid, and 6 for HalfCheetah. In P2P-RL, we adopt the
 31 normalization technique used in TD3+BC for the state and RL loss, and the hyperparameter α
 32 is set to 2. The model is trained 50 times for InvertedDoublePendulum, Hopper and HalfCheetah, and
 33 20 times for rest of the tasks. To learn the ζ and ν network in DualDICE, we first train them for
 34 $1e5$ times at the 5-th epoch and then train them for 2 times every time before policy learning. The
 35 learning rate of these two networks are both set to $1e-4$ and the batch size is set to 1024.

36 B.2 Environment Settings

37 All the environments remain the same with the original version of the tasks, except for the Inverted-
 38 DoublePendulum task where some additional noises are added to the states, which are set to Gaussian
 39 noises with mean 0 and standard deviation 10.

40 C Analysis of the Model Learning Process of MPC-RL

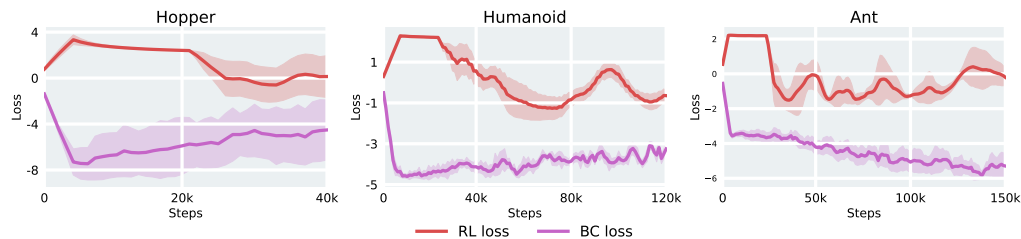


Figure 1: Quantitative analysis of the two kinds of losses in the model learning process of MPC-RL. “RL loss” means the loss of the reinforcement learning objective, and “BC loss” means the loss of the behavior cloning objective.

41 As shown in Section 4.1, the performance of P2P-RL is generally worse than P2P-MPC. An investiga-
 42 tion on the model learning process is shown in Figure 1, implying that the cause of this suboptimality
 43 may be the difficulty in balancing the loss of behavior cloning and RL.

44 D Extended Experiment for the Case When the Goal is in an Uncertain 45 Region

46 Serving as an extended experiment for the maze experiment in Section 4.4, we investigate the
 47 case when the goal is in an uncertain region in the online setting. For the convenience of im-

48 plementation, here the term "uncertainty" is equated with the epistemic uncertainty [Chua et al.,
 49 2018], which can be quantified as the amount of relevant real-world data. Therefore, a region
 50 with more data is considered to have lower uncertainty. Since in pure online settings the uncer-
 51 tainty of regions is hard to control during the training iterations, we first pretrain the model with
 52 an offline dataset and then switch to online training. The goal is allocated to the grey region
 53 in Figure 6 where the relevant offline samples are partially discarded. The percent of discarded
 54 samples is set to 25%, 50%, 75% and 100% respectively and the results are shown in Table 1.

	25%	50%	75%	100%
P2P-MPC	148.9 ± 35.9	75.4 ± 31.6	51.7 ± 29.8	43.2 ± 25.1
MBPO	116.2 ± 35.6	61.1 ± 34.8	47.5 ± 35.1	44.7 ± 30.2

Table 1: Results on the scenario where goal is in an uncertain region.

56 As the degree of uncertainty increases, the performances of both methods degrade rapidly, but P2P-
 57 MPC still outperforms MBPO in all these cases except for the 100% case, where P2P-MPC achieves
 58 slightly worse performance in average but better stability with lower standard deviation. To give a
 59 possible explanation of these results, it is worth noting that 1) P2P does not directly intervene the
 60 learning of policy or value function, but only improves the accuracy of the generated samples. As a
 61 result, the value function can still predict high value for uncertain regions and thus encourage the
 62 policy to explore them in the real environment; and 2) in contrast, even if the goal is in a region of
 63 high uncertainty and the model does not prevent the policy from exploring this region in the model,
 64 the value function can still predict low value of this region due to the lack of relevant data and thus
 65 mislead the learning of policy.

66 E Pseudo Code

67 The detailed descriptions of P2P-MPC and P2P-RL are respectively provided in Algorithm 1 and
 68 Algorithm 2.

Algorithm 1 P2P-MPC

```

1: Initialize policy  $\pi$ , predictive model  $\hat{P}$ , model-error predictor  $\hat{R}^m$ , environment dataset  $\mathcal{D}_e$  and
   model dataset  $\mathcal{D}_m$ .
2: for  $N$  epochs do
3:   Train model  $\hat{P}$  on  $\mathcal{D}_e$  via one-step prediction loss;
4:   Train  $\hat{R}^m$  on  $\mathcal{D}_e$ ;
5:   for  $E$  steps do
6:     Take action in environment according to  $\pi$ , and store the new transition to  $\mathcal{D}_e$ ;
7:     for  $M$  model rollouts do
8:       Sample initial states  $s \sim \mathcal{D}_e$  uniformly for rollout trajectories;
9:       for  $k$  rollout steps do
10:        Take action  $a$  according to  $\pi$  and the current state  $s$  in the model;
11:        Initialize  $s_0^m = (s, a)$  and perform  $L$  parallelized rollouts for  $H$  steps;
12:        Compute  $\sum_{t=0}^{H-1} \hat{R}^m(s_{t,j}^m, a_{t,j}^m)$  for each rollout trajectory  $j$ , denoted as  $r_j^m$ ,  $j \in$ 
            $\{1, 2, \dots, L\}$ , and take  $(s', r) = a_{0, \arg \max_j r_j^m}^m$ ;
13:        Store  $(s, a, r, s')$  to  $\mathcal{D}_m$  and then Let  $s = s'$ ;
14:       end for
15:     end for
16:     for  $G$  gradient updates do
17:       Update  $\pi$  using data sampled from  $\mathcal{D}_m$ ;
18:     end for
19:   end for
20: end for

```

Algorithm 2 P2P-RL

- 1: Initialize policy π , predictive model \hat{P} , environment dataset \mathcal{D}_e , model dataset \mathcal{D}_m , ζ and ν for using of DualDICE.
 - 2: **for** N epochs **do**
 - 3: Train the ζ and ν network according to the following objective derived by [Nachum et al., 2019]: $\mathbb{E}_{\mathcal{D}_e}[(\nu(s_t^m, a_t^m) - \gamma\nu(s_{t+1}^m, a_{t+1}^m))\zeta(s_t^m, a_t^m) - f^*(\zeta(s_t^m, a_t^m)) - (1 - \gamma)\nu(s_0^m, a_0^m)]$, where $f^*(x) := \frac{2}{3}|x|^{\frac{2}{3}}$ and s_{t+1}^m is updated from $(s_{t+1}, \pi_D(s_{t+1}))$ to $(s_{t+1}, \pi(s_{t+1}))$. Note that here π_D means the data-collecting policy and π the current policy;
 - 4: Train \hat{P} by optimizing $\mathbb{E}_{s_t^m, r_t^m \sim \mathcal{D}_e, a^m \sim \hat{P}(\cdot|s_t^m)}[\zeta(s_t^m, a^m)(\log \hat{P}(a^m|s_t^m) - Q^m(s_t^m, a^m)) + r_t^m]$, where the first term is the SAC [Haarnoja et al., 2018] loss and the second is the behavior cloning loss. Note that $Q^m(s^m, a^m)$ is the action-value function of \hat{P} and is trained by the same objective as the one in SAC;
 - 5: **for** E steps **do**
 - 6: Take action a_t in environment according to π and the current state s_t , then obtain the next state s_{t+1} and the reward r_{t+1} ;
 - 7: Reorder the transition: $(s_t, a_t, r_{t+1}, s_{t+1}, \hat{r}_{t+1}, \hat{s}_{t+1}) \rightarrow (s_t^m, a_t^m, r_{t+1}^m, s_{t+1}^m)$, where $s_t^m = (s_t, a_t)$, $a_t^m = (\hat{r}_{t+1}, \hat{s}_{t+1} - s_{t+1})$, $r_{t+1}^m = -\|\hat{s}_{t+1} - s_{t+1}\| - \|\hat{r}_{t+1} - r_{t+1}\|$, and $s_{t+1}^m = (s_{t+1}, a_{t+1})$ if $t + 1 \leq E$ else $(s_E, \pi(s_E))$; Store $(s_t^m, a_t^m, r_{t+1}^m, s_{t+1}^m)$ into \mathcal{D}_e ;
 - 8: **for** M model rollouts **do**
 - 9: Sample initial states uniformly from \mathcal{D}_e ;
 - 10: Perform k -step model rollouts starting from these states, using π and \hat{P} ; Add the generated samples to \mathcal{D}_m ;
 - 11: **end for**
 - 12: **for** G gradient updates **do**
 - 13: Update π using data sampled from \mathcal{D}_m ;
 - 14: **end for**
 - 15: **end for**
 - 16: **end for**
-

69 **References**

- 70 Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model: model-based
71 policy optimization. In *Proceedings of the 33rd International Conference on Neural Information*
72 *Processing Systems*, pages 12519–12530, 2019.
- 73 Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforcement
74 learning in a handful of trials using probabilistic dynamics models. In *Proceedings of the 32nd*
75 *International Conference on Neural Information Processing Systems*, pages 4759–4770, 2018.
- 76 Ofir Nachum, Yinlam Chow, Bo Dai, and Lihong Li. Dualdice: behavior-agnostic estimation of
77 discounted stationary distribution corrections. In *Proceedings of the 33rd International Conference*
78 *on Neural Information Processing Systems*, pages 2318–2328, 2019.
- 79 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
80 maximum entropy deep reinforcement learning with a stochastic actor. In *International conference*
81 *on machine learning*, pages 1861–1870. PMLR, 2018.