

382 Appendices

383 A Model-based Policy Optimization with Performance Guarantees

384 In this appendix, we provide proofs for bounds presented in the main paper.

385 We begin with a standard bound on model-based policy optimization, with bounded policy change ϵ_π
 386 and model error ϵ_m .

387 **Theorem A.1** (MBPO performance bound). *Let the expected total variation between two transition*
 388 *distributions is bounded at each timestep by $\max_t E_{s \sim \pi_{D,t}} [D_{TV}(p(s'|s, a) || \hat{p}(s'|s, a))] \leq \epsilon_m$, and*
 389 *the policy divergences are bounded as $\max_s D_{TV}(\pi_D(a|s) || \pi(a|s)) \leq \epsilon_\pi$. Then the returns are*
 390 *bounded as:*

$$\eta[\pi] \geq \hat{\eta}[\pi] - \frac{2\gamma r_{\max}(\epsilon_m + 2\epsilon_\pi)}{(1-\gamma)^2} - \frac{4r_{\max}\epsilon_\pi}{(1-\gamma)}$$

391 *Proof.* Let π_D denote the data collecting policy. As-is we can use Lemma B.3 to bound the returns,
 392 but it will require bounded model error under the new policy π . Thus, we need to introduce π_D by
 393 adding and subtracting $\eta[\pi_D]$, to get:

$$\eta[\pi] - \hat{\eta}[\pi] = \underbrace{\eta[\pi] - \eta[\pi_D]}_{L_1} + \underbrace{\eta[\pi_D] - \hat{\eta}[\pi]}_{L_2}$$

394 We can bound L_1 and L_2 both using Lemma B.3.

395 For L_1 , we apply Lemma B.3 using $\delta = \epsilon_\pi$ (no model error because both terms are under the true
 396 model), and obtain:

$$L_1 \geq -\frac{2r_{\max}\gamma\epsilon_\pi}{(1-\gamma)^2} - \frac{2r_{\max}\epsilon_\pi}{1-\gamma}$$

397 For L_2 , we apply Lemma B.3 using $\delta = \epsilon_\pi + \epsilon_m$ and obtain:

$$L_2 \geq -\frac{2r_{\max}\gamma(\epsilon_m + \epsilon_\pi)}{(1-\gamma)^2} - \frac{2r_{\max}\epsilon_\pi}{1-\gamma}$$

398 Adding these two bounds together yields the desired result. \square

399 Next, we describe bounds for branched rollouts. We define a branched rollout as a rollout which
 400 begins under some policy and dynamics (either true or learned), and at some point in time switches to
 401 rolling out under a new policy and dynamics for k steps. The point at which the branch is selected
 402 is weighted exponentially in time – that is, the probability of a branch point t being selected is
 403 proportional to γ^t .

404 We first present the simpler bound where the model error is bounded under the new policy, which we
 405 label as $\epsilon_{m'}$. This bound is difficult to apply in practice as supervised learning will typically control
 406 model error under the dataset collected by the previous policy.

407 **Theorem A.2.** *Let the expected total variation between two the learned model is bounded at each*
 408 *timestep under the expectation of π by $\max_t E_{s \sim \pi_t} [D_{TV}(p(s'|s, a) || \hat{p}(s'|s, a))] \leq \epsilon_{m'}$, and the*
 409 *policy divergences are bounded as $\max_s D_{TV}(\pi_D(a|s) || \pi(a|s)) \leq \epsilon_\pi$. Then under a branched*
 410 *rollouts scheme with a branch length of k , the returns are bounded as:*

$$\eta[\pi] \geq \eta^{\text{branch}}[\pi] - 2r_{\max} \left[\frac{\gamma^{k+1}\epsilon_\pi}{(1-\gamma)^2} + \frac{\gamma^k\epsilon_\pi}{(1-\gamma)} + \frac{k}{1-\gamma}(\epsilon_{m'}) \right]$$

411 *Proof.* As in the proof for Theorem A.1, the proof for this theorem requires adding and subtracting
 412 the correct reference quantity and applying the corresponding returns bound (Lemma B.4).

413 The choice of reference quantity is a branched rollout which executes the old policy π_D under the
 414 true dynamics until the branch point, then executes the new policy π under the true dynamics for k
 415 steps. We denote the returns under this scheme as $\eta^{\pi_D, \pi}$. We can split the returns as follows:

$$\eta[\pi] - \eta^{\text{branch}} = \underbrace{\eta[\pi] - \eta^{\pi_D, \pi}}_{L_1} + \underbrace{\eta^{\pi_D, \pi} - \eta^{\text{branch}}}_{L_2}$$

416 We can bound both terms L_1 and L_2 using Lemma B.4.

417 L_1 accounts for the error from executing the old policy instead of the current policy. This term only
 418 suffers from error *before* the branch begins, and we can use Lemma B.4 $\epsilon_\pi^{\text{pre}} \leq \epsilon_\pi$ and all other errors
 419 set to 0. This implies:

$$|\eta[\pi] - \eta^{\pi_D, \pi}| \leq 2r_{\max} \left[\frac{\gamma^{k+1}}{(1-\gamma)^2} \epsilon_\pi + \frac{\gamma^k}{1-\gamma} \epsilon_\pi \right]$$

420 L_2 incorporates model error *under the new policy* incurred after the branch. Again we use Lemma B.4,
 421 setting $\epsilon_m^{\text{post}} \leq \epsilon_m$ and all other errors set to 0. This implies:

$$|\eta[\pi] - \eta^{\pi_D, \pi}| \leq 2r_{\max} \left[\frac{k}{1-\gamma} \epsilon_{m'} \right]$$

422 Adding L_1 and L_2 together completes the proof.

423 □

424 The next bound is an analogue of Theorem A.2 except using model errors under the previous policy
 425 π_D rather than the new policy π .

426 **Theorem A.3.** *Let the expected total variation between two the learned model is bounded at each*
 427 *timestep under the expectation of π by $\max_t E_{s \sim \pi_D, t} [D_{TV}(p(s'|s, a) || \hat{p}(s'|s, a))] \leq \epsilon_m$, and the*
 428 *policy divergences are bounded as $\max_s D_{TV}(\pi_D(a|s) || \pi(a|s)) \leq \epsilon_\pi$. Then under a branched*
 429 *rollouts scheme with a branch length of k , the returns are bounded as:*

$$\eta[\pi] \geq \eta^{\text{branch}}[\pi] - 2r_{\max} \left[\frac{\gamma^{k+1} \epsilon_\pi}{(1-\gamma)^2} + \frac{\gamma^k + 2}{(1-\gamma)} \epsilon_\pi + \frac{k}{1-\gamma} (\epsilon_m + 2\epsilon_\pi) \right]$$

430 *Proof.* This proof is a short extension of the proof for Theorem A.2. The only modification is that
 431 we need to bound L_2 in terms of the model error under the π_D rather than π .

432 Once again, we design a new reference rollout. We use a rollout that executes the old policy π_D
 433 under the true dynamics until the branch point, then executes the *old* policy π_D under the model for
 434 k steps. We denote the returns under this scheme as $\eta^{\pi_D, \hat{\pi}_D}$. We can split L_2 as follows:

$$\eta^{\pi_D, \pi} - \eta^{\text{branch}} = \underbrace{\eta^{\pi_D, \pi} - \eta^{\pi_D, \hat{\pi}_D}}_{L_3} + \underbrace{\eta^{\pi_D, \hat{\pi}_D} - \eta^{\text{branch}}}_{L_4}$$

435 Once again, we bound both terms L_3 and L_4 using Lemma B.4.

436 The rollouts in L_3 differ in both model and policy after the branch. This can be bound using
 437 Lemma B.4 by setting $\epsilon_\pi^{\text{post}} = \epsilon_\pi$ and $\epsilon_m^{\text{post}} = \epsilon_m$. This results in:

$$|\eta^{\pi_D, \pi} - \eta^{\pi_D, \hat{\pi}_D}| \leq 2r_{\max} \left[\frac{k}{1-\gamma} (\epsilon_m + \epsilon_\pi) + \frac{1}{1-\gamma} \epsilon_\pi \right]$$

438 The rollouts in L_4 differ only in the policy after the branch (as they both rollout under the model).
 439 This can be bound using Lemma B.4 by setting $\epsilon_\pi^{\text{post}} = \epsilon_\pi$ and $\epsilon_m^{\text{post}} = 0$. This results in:

$$|\eta^{\pi_D, \hat{\pi}_D} - \eta^{\text{branch}}| \leq 2r_{\max} \left[\frac{k}{1-\gamma} (\epsilon_\pi) + \frac{1}{1-\gamma} \epsilon_\pi \right]$$

440 Adding L_1 from Theorem A.2 and L_3, L_4 above completes the proof.

441 □

442 **B Useful Lemmas**

443 In this section, we provide proofs for various lemmas used in our bounds.

444 **Lemma B.1** (TVD of Joint Distributions). *Suppose we have two distributions $p_1(x, y) =$
 445 $p_1(x)p_1(y|x)$ and $p_2(x, y) = p_2(x)p_2(y|x)$. We can bound the total variation distance of the
 446 joint as:*

$$D_{TV}(p_1(x, y)||p_2(x, y)) \leq D_{TV}(p_1(x)||p_2(x)) + \max_x D_{TV}(p_1(y|x)||p_2(y|x))$$

447 *Alternatively, we have a tighter bound in terms of the expected TVD of the conditional:*

$$D_{TV}(p_1(x, y)||p_2(x, y)) \leq D_{TV}(p_1(x)||p_2(x)) + E_{x \sim p_1}[D_{TV}(p_1(y|x)||p_2(y|x))]$$

Proof.

$$\begin{aligned} D_{TV}(p_1(x, y)||p_2(x, y)) &= \frac{1}{2} \sum_{x, y} |p_1(x, y) - p_2(x, y)| \\ &= \frac{1}{2} \sum_{x, y} |p_1(x)p_1(y|x) - p_2(x)p_2(y|x)| \\ &= \frac{1}{2} \sum_{x, y} |p_1(x)p_1(y|x) - p_1(x)p_2(y|x) + (p_1(x) - p_2(x))p_2(y|x)| \\ &\leq \frac{1}{2} \sum_{x, y} p_1(x)|p_1(y|x) - p_2(y|x)| + |p_1(x) - p_2(x)|p_2(y|x) \\ &= \frac{1}{2} \sum_{x, y} p_1(x)|p_1(y|x) - p_2(y|x)| + \frac{1}{2} \sum_x |p_1(x) - p_2(x)| \\ &= E_{x \sim p_1}[D_{TV}(p_1(y|x)||p_2(y|x))] + D_{TV}(p_1(x)||p_2(x)) \\ &\leq \max_x D_{TV}(p_1(y|x)||p_2(y|x)) + D_{TV}(p_1(x)||p_2(x)) \end{aligned}$$

448 □

449 **Lemma B.2** (Markov chain TVD bound, time-varying). *Suppose the expected KL-divergence between
 450 two transition distributions is bounded as $\max_t E_{s \sim p_1^t(s)} D_{KL}(p_1(s'|s)||p_2(s'|s)) \leq \delta$, and the
 451 initial state distributions are the same – $p_1^{t=0}(s) = p_2^{t=0}(s)$. Then the distance in the state marginal
 452 is bounded as:*

$$D_{TV}(p_1^t(s)||p_2^t(s)) \leq t\delta$$

453 *Proof.* We begin by bounding the TVD in state-visitation at time t , which is denoted as $\epsilon_t =$
 454 $D_{TV}(p_1^t(s)||p_2^t(s))$.

$$\begin{aligned} |p_1^t(s) - p_2^t(s)| &= \left| \sum_{s'} p_1(s_t = s|s')p_1^{t-1}(s') - p_2(s_t = s|s')p_2^{t-1}(s') \right| \\ &\leq \sum_{s'} |p_1(s_t = s|s')p_1^{t-1}(s') - p_2(s_t = s|s')p_2^{t-1}(s')| \\ &= \sum_{s'} |p_1(s|s')p_1^{t-1}(s') - p_2(s|s')p_1^{t-1}(s') + p_2(s|s')p_1^{t-1}(s') - p_2(s|s')p_2^{t-1}(s')| \\ &\leq \sum_{s'} p_1^{t-1}(s')|p_1(s|s') - p_2(s|s')| + p_2(s|s')|p_1^{t-1}(s') - p_2^{t-1}(s')| \\ &= E_{s' \sim p_1^{t-1}}[|p_1(s|s') - p_2(s|s')|] + \sum_{s'} p(s|s')|p_1^{t-1}(s') - p_2^{t-1}(s')| \end{aligned}$$

$$\begin{aligned}
\epsilon_t &= D_{TV}(p_1^t(s)||p_2^t(s)) = \frac{1}{2} \sum_s |p_1^t(s) - p_2^t(s)| \\
&= \frac{1}{2} \sum_s \left(E_{s' \sim p_1^{t-1}} [|p_1(s|s') - p_2(s|s')|] + \sum_{s'} p(s|s') |p_1^{t-1}(s') - p_2^{t-1}(s')| \right) \\
&= \frac{1}{2} E_{s' \sim p_1^{t-1}} \left[\sum_s |p_1(s|s') - p_2(s|s')| \right] + D_{TV}(p_1^{t-1}(s')||p_2^{t-1}(s')) \\
&= \delta_t + \epsilon_{t-1} \\
&= \epsilon_0 + \sum_{i=0}^t \delta_i \\
&= \sum_{i=0}^t \delta_i = t\delta
\end{aligned}$$

456 Where we have defined $\delta_t = \frac{1}{2} E_{s' \sim p_1^{t-1}} [\sum_s |p_1(s|s') - p_2(s|s')|]$, which we assume is upper bounded
457 by δ . Assuming we are not modeling the initial state distribution, we can set $\epsilon_0 = 0$. \square

458 **Lemma B.3** (Branched Returns bound). *Suppose the expected KL-divergence between two*
459 *dynamics distributions is bounded as $\max_t E_{s \sim p_1^t(s)} D_{KL}(p_1(s', a|s)||p_2(s', a|s)) \leq \epsilon_m$, and*
460 *$\max_s D_{TV}(\pi_1(a|s)||\pi_2(a|s)) \leq \epsilon_\pi$. Then the returns are bounded as:*

$$|\eta_1 - \eta_2| \leq \frac{2R\gamma(\epsilon_\pi + \epsilon_m)}{(1 - \gamma)^2} + \frac{2R\epsilon_\pi}{1 - \gamma}$$

461 *Proof.* Here, η_1 denotes returns of π_1 under dynamics $p_1(s'|s, a)$, and η_2 denotes returns of π_2 under
462 dynamics $p_2(s'|s, a)$.

$$\begin{aligned}
|\eta_1 - \eta_2| &= \left| \sum_{s,a} (p_1(s, a) - p_2(s, a)) r(s, a) \right| \\
&= \left| \sum_{s,a} \left(\sum_t \gamma^t p_1^t(s, a) - p_2^t(s, a) \right) r(s, a) \right| \\
&= \left| \sum_t \sum_{s,a} \gamma^t (p_1^t(s, a) - p_2^t(s, a)) r(s, a) \right| \\
&\leq \sum_t \sum_{s,a} \gamma^t |p_1^t(s, a) - p_2^t(s, a)| r(s, a) \\
&\leq r_{\max} \sum_t \sum_{s,a} \gamma^t |p_1^t(s, a) - p_2^t(s, a)|
\end{aligned}$$

463 We now apply Lemma B.2, using $\delta = \epsilon_m + \epsilon_\pi$ (via Lemma B.1) to get:

$$D_{TV}(p_1^t(s)||p_2^t(s)) \leq t(\epsilon_m + \epsilon_\pi)$$

464 And since we assume $\max_s D_{TV}(\pi_1(a|s)||\pi_2(a|s)) \leq \epsilon_\pi$, we get

$$D_{TV}(p_1^t(s, a)||p_2^t(s, a)) \leq t(\epsilon_m + \epsilon_\pi) + \epsilon_\pi$$

465 Thus, plugging this back in we get:

$$\begin{aligned}
|\eta_1 - \eta_2| &\leq r_{\max} \sum_t \sum_{s,a} \gamma^t |p_1^t(s, a) - p_2^t(s, a)| \\
&\leq 2r_{\max} \sum_t \gamma^t t(\epsilon_m + \epsilon_\pi) + \epsilon_\pi \\
&\leq 2r_{\max} \left(\frac{\gamma(\epsilon_\pi + \epsilon_m)}{(1 - \gamma)^2} + \frac{\epsilon_\pi}{1 - \gamma} \right)
\end{aligned}$$

467 **Lemma B.4** (Returns bound, branched rollout). Assume we run a branched rollout of length
468 k . Before the branch (“pre” branch), we assume that the dynamics distributions are
469 bounded as $\max_t E_{s \sim p_1^t(s)} D_{KL}(p_1^{\text{pre}}(s', a|s) || p_2^{\text{pre}}(s', a|s)) \leq \epsilon_m^{\text{pre}}$ and after the branch as
470 $\max_t E_{s \sim p_1^t(s)} D_{KL}(p_1^{\text{post}}(s', a|s) || p_2^{\text{post}}(s', a|s)) \leq \epsilon_m^{\text{post}}$. Likewise, the policy divergence is
471 bounded pre- and post- branch by $\epsilon_\pi^{\text{pre}}$ and $\epsilon_\pi^{\text{post}}$, respectively. Then the K -step returns are bounded
472 as:

$$|\eta_1 - \eta_2| \leq 2r_{\max} \left[\frac{\gamma^{k+1}}{(1-\gamma)^2} (\epsilon_m^{\text{pre}} + \epsilon_\pi^{\text{pre}}) + \frac{k}{1-\gamma} (\epsilon_m^{\text{post}} + \epsilon_\pi^{\text{post}}) + \frac{\gamma^k}{1-\gamma} \epsilon_\pi^{\text{pre}} + \frac{1}{1-\gamma} \epsilon_\pi^{\text{post}} \right]$$

473 *Proof.* We begin by bounding state marginals at each timestep, similar to Lemma B.3. Recall that
474 Lemma B.2 implies that state marginal error at each timestep can be bounded by the state marginal
475 error at the previous timestep, plus the divergence at the current timestep.

476 Thus, letting $d_1(s, a)$ and $d_2(s, a)$ denote the state-action marginals, we can write:

477 For $t \leq k$:

$$TV d_1^t(s, a) d_2^t(s, a) \leq t(\epsilon_m^{\text{post}} + \epsilon_\pi^{\text{post}}) + \epsilon_\pi^{\text{post}} \leq k(\epsilon_m^{\text{post}} + \epsilon_\pi^{\text{post}}) + \epsilon_\pi^{\text{post}}$$

478 and for $t \geq k$:

$$TV d_1^t(s, a) d_2^t(s, a) \leq (t-k)(\epsilon_m^{\text{pre}} + \epsilon_\pi^{\text{pre}}) + k(\epsilon_m^{\text{post}} + \epsilon_\pi^{\text{post}}) + \epsilon_\pi^{\text{pre}} + \epsilon_\pi^{\text{post}}$$

479 We can now bound the difference in occupancy measures by averaging the state marginal error over
480 time, weighted by the discount:

$$\begin{aligned} D_{TV}(d_1(s, a) || d_2(s, a)) &\leq (1-\gamma) \sum_{t=0}^{\infty} \gamma^t t D_{TV}(d_1^t(s, a) || d_2^t(s, a)) \\ &\leq (1-\gamma) \sum_{t=0}^k \gamma^t (k(\epsilon_m^{\text{post}} + \epsilon_\pi^{\text{post}}) + \epsilon_\pi^{\text{post}}) \\ &\quad + (1-\gamma) \sum_{t=k}^{\infty} \gamma^t (t-k)(\epsilon_m^{\text{pre}} + \epsilon_\pi^{\text{pre}}) + k(\epsilon_m^{\text{post}} + \epsilon_\pi^{\text{post}}) + \epsilon_\pi^{\text{pre}} + \epsilon_\pi^{\text{post}} \\ &= k(\epsilon_m^{\text{post}} + \epsilon_\pi^{\text{post}} + \epsilon_\pi^{\text{post}}) + \frac{\gamma^{k+1}}{1-\gamma} (\epsilon_m^{\text{pre}} + \epsilon_\pi^{\text{pre}}) + \gamma^k \epsilon_\pi^{\text{pre}} \end{aligned}$$

481 Multiplying this bound by $\frac{2r_{\max}}{1-\gamma}$ to convert the state-marginal bound into a returns bound completes
482 the proof. \square

483 **C Hyperparameter Settings**

		HalfCheetah	Walker2d	Ant	Hopper
N	epochs	400		300	125
E	environment steps per epoch	1000			
M	model rollouts per policy update	20			
B	ensemble size	7			
	network architecture	MLP with four hidden layers of size 200			
G	policy updates per epoch	40000		20000	
k	model horizon	1		1 \rightarrow 25 over epochs 20 \rightarrow 100	1 \rightarrow 15 over epochs 20 \rightarrow 100

Table 1: Hyperparameter settings for MBPO results shown in Figure 2. $x \rightarrow y$ over epochs $a \rightarrow b$ denotes a thresholded linear function, *i.e.* at epoch e , $f(e) = \min(\max(x + \frac{e-a}{b-a} \cdot (x-y), x), y)$