
Discriminator Augmented Model-Based Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 By planning through a learned dynamics model, model-based reinforcement learning
2 (MBRL) offers the prospect of good performance with little environment
3 interaction. However, it is common in practice for the learned model to be in-
4 accurate, impairing planning and leading to poor performance. This paper aims
5 to improve planning with an importance sampling framework that accounts and
6 corrects for discrepancy between the true and learned dynamics. This framework
7 also motivates an alternative objective for fitting the dynamics model: to minimize
8 the variance of value estimation during planning. We derive and implement this
9 objective, which encourages better prediction on trajectories with larger returns.
10 We observe empirically that our approach improves the performance of current
11 MBRL algorithms on two stochastic control problems, and provide a theoretical
12 basis for our method.

13 1 Introduction

14 Model free reinforcement learning methods have achieved good performance on a number of complex
15 tasks, but usually require a large amount of data collected through environment interaction [25, 23].
16 *Model-based* reinforcement learning (MBRL) can potentially reduce these data requirements by fitting
17 a model of the environment dynamics on a small dataset, then planning through the learned dynamics
18 to produce a good policy. However, inaccurate models can severely impair planning performance.
19 Recent techniques attempt to avoid *compounding* model error by restricting the number of model
20 unrolling steps [16, 11], but this creates a trade-off between planning performance and sample
21 efficiency. More importantly, such approaches assume that the learned model mostly matches the true
22 (one timestep) dynamics, and only address compounding error that arises over many timesteps. In this
23 paper we show that there are a range of environments where the learned model can be inaccurate even
24 on short time horizons (e.g., one timestep), so merely mitigating compounding error is insufficient.
25 Given the evidence that even deep neural network models can struggle to learn complex, high
26 dimensional distributions [2], we expect inaccurate learned dynamics to be an increasingly important
27 issue as we apply MBRL to harder and more realistic environments.

28 Consider a learned dynamics model as a generative model from which we sample transitions. Even if
29 the model is inaccurate, some samples may be better than others. Our approach trains a discriminative
30 model to assess the quality of sampled transitions during planning, and upweight or downweight value
31 estimates computed from high and low quality samples, respectively. In section 3.1 we also show
32 that this method is a form of likelihood-free importance sampling that, assuming the discriminator
33 is an optimal classifier, produces unbiased value estimates for planning. Since our discriminator
34 can correct model error during planning, we are no longer restricted to fitting the dynamics using a
35 maximum likelihood objective. Instead, we can learn *biased* dynamics models with advantageous
36 properties, such as reduced value estimation variance during planning. We derive and implement an

37 objective function for learning these variance-minimizing models. Unlike maximum likelihood, our
 38 proposed objective incorporates trajectory return statistics into the model fitting process. We call our
 39 planning and model training framework **Discriminator Augmented MBRL** (or **DAM** for short).

40 To evaluate scenarios where the dynamics is difficult to learn, our experiments consider environments
 41 with stochastic and multi-modal dynamics. We find that DAM significantly improves planning
 42 performance over existing MBRL algorithms in these environments.

43 2 Preliminaries

44 Consider a Markov Decision Process (MDP) $\mathcal{M} \equiv (\mathcal{S}, \mathcal{A}, p, r, \gamma, T)$ [30], where \mathcal{S} denotes the state
 45 space, \mathcal{A} denotes the action space, $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto \mathbb{R}_{>0}$ denotes the transition dynamics, $r : \mathcal{S} \times \mathcal{A} \mapsto$
 46 \mathbb{R} denotes the reward function, $\gamma \in [0, 1]$ denotes the discount factor and T denotes the maximum
 47 episode horizon. Our objective is to maximize $\mathbb{E}_{p(\tau)}[R(\tau)]$ where $R(\tau) = \sum_{t=1}^T \gamma^t r(s_t, a_t)$ denotes
 48 the discounted return for the trajectory $\tau = \{s_1, a_1, \dots, s_T\}$, which is sampled from

$$p(\tau) = p(s_1) \prod_{t=1}^{T-1} p(s_{t+1}|s_t, a_t) \pi(a_t|s_t). \quad (1)$$

49 Here, $p : \mathcal{S} \mapsto \mathbb{R}_{>0}$ denotes the initial state distribution, and $\pi(a_t | s_t)$ defines our controller which
 50 is optimized to maximize the expected discounted return. The overloaded notation p is interpreted
 51 based on variable(s) for which the density is computed.

52 Since p can only be sampled sequentially in practice, MBRL [20, 28] learns an estimate of the transi-
 53 tion dynamics $q : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto \mathbb{R}_{>0}$. We can also define $q(\tau) = p(s_1) \prod_{t=1}^{T-1} q(s_{t+1}|s_t, a_t) \pi(a_t|s_t)$
 54 analogous to $p(\tau)$. The model is generally learned by maximizing $\mathbb{E}_{s', a, s \sim \mathcal{B}} [\log q(s' | s, a)]$, where
 55 \mathcal{B} denotes the dataset of transitions collected in \mathcal{M} . The approximate transition dynamics can be
 56 used to learn a parametric controller π or perform online planning using sampling based methods.

57 3 Discriminator Augmented MBRL

Function ESTIMATEVALUE

```

input :  $s_1$ : Initial state
input :  $\mathbf{a} = \{a_t\}_{t=1}^H$ : Action sequence
 $w \leftarrow 1$ ; // stores weights (Eq. 4)
 $R \leftarrow 0$ ;
for  $t \leftarrow 1$  to  $H$  do
   $r_t \leftarrow r(s_t, a_t)$ ;
   $R \leftarrow R + w \cdot r_t$ ;
   $s_{t+1} \sim q(\cdot | s_t, a_t)$ ;
   $w \leftarrow w \cdot \sigma(\mathcal{D}(s_t, a_t, s_{t+1}))$ ; // Eq. 8
end
Result:  $R$ 

```

Algorithm 1: ESTIMATEVALUE produces an estimated return R of some action sequence using a learned model q . Since q may be biased with respect to the true dynamics, it uses a discriminator \mathcal{D} to correct the estimate as described in Sec. 3.1

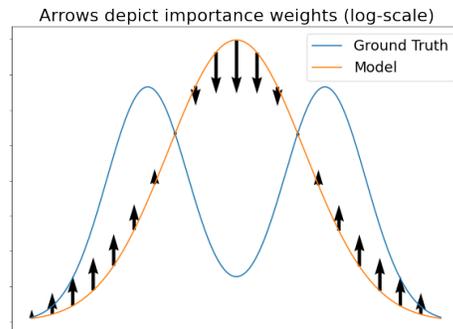


Figure 1: In this illustration, the blue curve depicts a ground truth density and orange a learned model. Importance sampling weights (represented by each arrow’s magnitude in log-scale) can help us correct estimates under the incorrect model distribution.

59 In this work, we look at some of the fundamental components of model-based reinforcement learning.
 60 In Section 3.1 we discuss an often ignored discrepancy caused by sampling trajectories from learned
 61 dynamics models and discuss how this bias can be corrected using importance sampling. Then, in
 62 Section 3.2 we discuss how interpreting the model-based RL in an importance sampling framework
 63 suggests a novel objective for learning models for environment dynamics.

64 **3.1 Correcting the Sampling Bias in Model-Based RL**

65 Consider our optimization objective $J(\pi) = \mathbb{E}_{p(\tau)}[R(\tau)]$. Using the approximate dynamics q intro-
 66 duces a bias in our estimation of $J(\pi)$. Model-based RL methods generally use $\tilde{J}_q(\pi) = \mathbb{E}_{q(\tau)}[R(\tau)]$
 67 as a proxy for $J(\pi)$, where the trajectories are sampled from $q(\tau)$ instead of $p(\tau)$ without accounting
 68 for the resulting bias. In this work, we mitigate this bias by using the importance sampling framework:

$$\mathbb{E}_{p(\tau)}[R(\tau)] = \mathbb{E}_{q(\tau)}\left[\frac{p(\tau)}{q(\tau)}R(\tau)\right] = \mathbb{E}_{q(\tau)}[w(\tau)R(\tau)] \quad (2)$$

69 where we have introduced the importance sampling correction $w(\tau)$ to correct for the bias introduced
 70 by sampling trajectories from $q(\tau)$. Simplifying $w(\tau)$, we get

$$w(\tau) = \prod_{t=1}^{T-1} \frac{p(s_{t+1} | s_t, a_t)}{q(s_{t+1} | s_t, a_t)} = \prod_{t=1}^{T-1} w(s_{t+1}, a_t, s_t) \quad (3)$$

71 By exploiting the MDP’s temporal structure, we can obtain an improved importance sampling
 72 estimator that weights the per-timestep rewards instead of the trajectory return. Comparing to Eq. 2,
 73 we see that this estimator multiplies the reward at time t by a weight that only depends on transitions
 74 leading up to it, and not transitions that occur after it:

$$\mathbb{E}_{p(\tau)}[R(\tau)] = \sum_{t=1}^T \mathbb{E}_q \left[\left(\prod_{m=1}^t w(s_m, a_m, s_{m+1}) \right) r(s_t, a_t) \right] \quad (4)$$

75 As stated earlier, we only have access to the samples from the transition dynamics and not the
 76 probabilities. Thus to compute $w(\tau)$ for $\tau \sim q$, we use techniques from density ratio estimation [36,
 77 13, 33]. We setup a binary classification problem over $\mathcal{D} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \mathcal{Y} \mapsto [0, 1]$, where
 78 $\mathcal{Y} = \{0, 1\}$. Using the transition samples from the environment and the learned dynamics, we create
 79 a dataset of transitions $\{s', a, s, y\}$, where $y = 1$ for $(s', a, s) \sim p(s' | s, a)p(s, a)$ and $y = 0$ for
 80 $(s', a, s) \sim q(s' | s, a)p(s, a)$. Effectively, we are training a discriminator to distinguish between real
 81 transitions in the environment and the transitions sampled from the learned model. Alg. 2 provides
 82 the pseudocode for training the discriminator.

83 To connect the discriminator to the importance weights, consider the Bayes-optimal classifier for this
 84 problem:

$$\mathcal{D}(s', a, s) = \frac{p(s' | s, a)p(s, a)}{p(s' | s, a)p(s, a) + q(s' | s, a)p(s, a)} \quad (5)$$

$$= \frac{p(s' | s, a)/q(s' | s, a)}{p(s' | s, a)/q(s' | s, a) + 1} \quad (6)$$

$$= \sigma \left(\log \frac{p(s' | s, a)}{q(s' | s, a)} \right) = \sigma \left(\tilde{\mathcal{D}}(s', a, s) \right) \quad (7)$$

85 $\tilde{\mathcal{D}}$ is the classifier logits, and we can recover the importance weight by exponentiating them:

$$w(s', a, s) = \frac{p(s', a, s)}{q(s', a, s)} = \exp(\tilde{\mathcal{D}}(s', a, s)) \quad (8)$$

86 For our transitions $\{s', a, s\}$, it is important to sample $s, a \sim p(s, a)$ (that is, real state-action pairs)
 87 so that the logits can be connected to the importance weights. Since, we only need samples to train
 88 the classifier, this allows us to estimate the importance weights without requiring probabilities from
 89 underlying transition dynamics. Alg. 1 provides the pseudocode for estimating the discriminator-
 90 corrected return of an action sequence, which forms the subroutine for choosing actions in sampling
 91 based planning algorithms.

92 The benefit of this approach arises from the potential simplicity of the discriminative problem
 93 compared to the modelling the dynamics. The dynamics can be hard to model in the real world,
 94 as they often involve higher dimensional spaces (like RGB images) and can be multi-modal due to
 95 partial observability. The modelling and optimization challenges can often result in an incorrectly
 96 learned model, an example of which is shown in Fig 1. On the other hand, discrimination can be

97 an easier problem as fake and real samples can be well separated, especially in higher-dimensional
 98 spaces. This can improve the value estimation by down-weighting samples which do not fit with the
 99 real world samples. If the learned dynamics are indeed correct, the true and fake samples will be
 100 indistinguishable and the importance sampling weights will be close to unity, reducing the problem
 101 to conventional model-based RL.

102 3.2 Minimum Variance Dynamics Model Fitting

103 The above discussion applies to any learned model, including the conventional approach of learning
 104 by maximizing the log-likelihood of transitions collected in the environment. The policy objective
 105 under q , that is $\hat{J}_q(\pi) = \frac{1}{N} \sum_{i=1}^N w(\tau_i) R(\tau_i)$, is an unbiased estimator of $J(\pi) \forall q$ as long as the
 106 support of p is a subset of support of q . However, from an optimization perspective, it is desirable
 107 to have an estimator with the minimum variance. Sampling based planning algorithms generate
 108 candidate action sequences and rank them according to their value function. Minimizing the variance
 109 of the value estimate of an action sequence increases the likelihood of the action sequences getting
 110 ranked correctly, and thus, the best action sequence getting chosen. Thus, it is favorable to reduce the
 111 variance of our estimate \hat{J}_q . Assuming that the $R(\tau) \geq 0$, it can be shown that

$$q^*(\tau) = \frac{R(\tau)p(\tau)}{\mathbb{E}_{p(\tau)}[R(\tau)]} \quad (9)$$

112 Here, q^* denotes the sampling distribution which minimizes the variance of \hat{J}_q . While the denominator
 113 for the optimal sampling distribution is intractable, the important distinction here is that trajectories
 114 should be sampled in accordance to their density under the true trajectory distribution and the return
 115 accumulated by the trajectory. Using the q^* as the target distribution, we can setup our minimization
 116 objective for model-based learning to be $\text{KL}(q^*(\tau) \parallel q(\tau))$. As we show in Appendix [A.1](#), the
 117 gradient with respect to q is:

$$\nabla_q \text{KL}(q^*(\tau) \parallel q(\tau)) \propto - \sum_{t=1}^{T-1} \mathbb{E}_{p(\tau)} [R(\tau) \nabla_q \log q(s_{t+1} \mid s_t, a_t)] \quad (10)$$

118 Intuitively, this corresponds to upweighting the gradient for trajectories with higher return, forcing
 119 the dynamics model to fit better to transitions with higher return. In particular, this highlights that
 120 the conventional approach of learning dynamics does not represent the optimal approach in the
 121 importance sampling framework. The modified model training algorithm is described in Alg. [2](#) in the
 122 appendix.

123 3.3 Planning and Training Loop

124 The previous two sections describe how to train the dynamics model, and how to train a discriminator
 125 that corrects the sampling bias from that learned model. Further, ESTIMATEVALUE (Alg. [1](#)) shows
 126 how to implement discriminator-corrected value estimation. We can plug ESTIMATEVALUE as a
 127 subroutine into most sampling based planners. Once a planner is in place, we can use it execute
 128 actions in the real environment, collect more data, and re-train our model and discriminator. In
 129 Appendix [B](#) we discuss in detail the implementation of this model-based RL loop within our model
 130 and discriminator training framework.

131 4 Experiments

132 In this section we evaluate our proposed approach and compare with model-based baselines. We
 133 hypothesize that environments with multi-modal dynamics make model fitting difficult, which would
 134 pose problems for planning. Indeed, we observe that our proposed approach improves performance
 135 over standard model-based RL in such settings.

136 We devise and implement two environments with stochastic, multi-modal dynamics: **IcyRoad** and
 137 **Intersection**. In **IcyRoad** the agent drives a car along an icy road, and above a threshold speed the car
 138 has a probability of swerving off the road in either direction (hence the multi-modality) and incurring
 139 a large negative reward. In **Intersection** the agent drives a car into an intersection. An oncoming car
 140 is also entering the intersection and will turn either left or right with equal probability—the outcome is

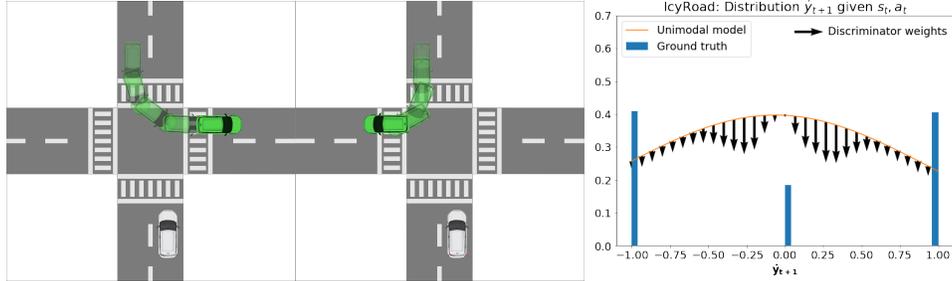


Figure 2: **Left:** Depiction of different scenarios in the Intersection environment, where the agent controls the silver car. Both cars are entering the intersection at the same time, but the agent does not know whether the oncoming green car will turn left (left panel) or turn right (right panel). If the agent drives into the intersection too fast, it may collide. **Right:** The learned (unimodal) model and ground truth distributions $p(\hat{y}_{s+1}|s_t, a_t)$, where $s_t = (2, 0, 2, 0)$ and the action is to accelerate. After training with maximum likelihood, the unimodal model erroneously places large probability mass on unrealistic regions of the state space. The discriminator produces importance weights (black arrows, log scale) that downweight the model’s samples in the unrealistic regions.

Average return	IcyRoad		Intersection	
	Unimodal	Multimodal	Unimodal	Multimodal
Baseline	0.71 ± 1.72	5.34 ± 0.84	62.67 ± 22.10	29.58 ± 0.105
DAM (minvar only)	0.86 ± 0.79	4.55 ± 3.50	31.40 ± 1.111	31.30 ± 1.634
DAM (disc only)	4.88 ± 0.50	3.30 ± 2.75	87.79 ± 15.34	60.18 ± 19.58
DAM	2.32 ± 1.76	7.94 ± 0.92	46.76 ± 9.543	97.72 ± 5.127

Table 1: Average return for each method in the IcyRoad and Intersection environments. Columns are split by whether they use a unimodal (Gaussian) or multimodal (mixture density network) learned dynamics model. Combining our method **DAM** with multimodal learned dynamics models usually performs best among different methods we evaluated. Results show mean and standard error across 3 runs (random seeds) for each hyperparameter setting.

141 unknown in advance to the agent. Fig. 2 illustrates the Intersection environment, and Appendix C
 142 both environments in more detail.

143 For our baseline, we fit a dynamics model using the usual maximum likelihood loss and use a
 144 standard random shooting planner with uncorrected value estimates. In **DAM** we use both aspects of
 145 our method: we train the model using the minimum variance objective (Eq. 10) and plan using the
 146 discriminator-corrected estimator (Eq. 4). To disentangle the effects of each component of our method,
 147 we also evaluate variants that use only one aspect of our method, but not both. **DAM (disc only)** only
 148 does discriminator-corrected planning, while **DAM (minvar only)** only uses the minimum variance
 149 object. Practitioners typically train models that output the parameters for a Gaussian distribution,
 150 which is unimodal. Since our environment dynamics are multimodal, it also makes sense to try fitting
 151 multimodal models such as mixture density networks [4]. We test each of our planning and training
 152 methods with both a mixture density network model and a standard Gaussian model, and display both
 153 results.

154 Table 1 depicts the results after 10 rounds of interleaving model (+ discriminator) training with
 155 data collection in the environment. We see that in both environments, **DAM** using multimodal
 156 dynamics models obtains the best results. **DAM (minvar only)**, which only uses the minimum
 157 variance objective without discriminator-corrected planning, is often *worse* than the baseline. We
 158 hypothesize this occurs because the minimum variance objective increases model bias versus standard
 159 maximum likelihood training, which can be harmful without discriminator-corrected planning. When
 160 using unimodal models in either environment, we see that simply having a discriminator for planning
 161 already increases performance over the baseline. Fig. 2 depicts our analysis of what the trained model
 162 and discriminator are doing in IcyRoad: we see that the unimodal model is struggling to fit the true
 163 dynamics, while the discriminator is helping to correct the model error. Interestingly, adding the
 164 minimum variance objective to discriminator-corrected planning does *not* help with unimodal models,
 165 but it *does* help with multimodal models.

166 5 Related Work

167 Model-based RL has a rich history in control and robotics. Classically, the methods for model-based
168 control have assumed access to underlying dynamics [35, 22, 26, 24] or make simplifying assumptions
169 about the environment dynamics [7, 8, 21, 18]. While such simplifying assumptions can yield exact or
170 efficient controllers, the models cannot faithfully represent the underlying dynamics, especially when
171 considering high dimensional state spaces like RGB images. Recent work has focused on combining
172 high-capacity function approximators such as neural networks with model-based reinforcement
173 learning [20, 27, 28, 6, 12, 32, 14, 29], which employ sampling-based planning methods [38, 27, 31].
174 Concurrently, some recent work has integrated policy networks with deep models [16, 37] in spirit
175 of Dyna [34]. Neural networks allow more flexible representation of non-linear dynamics, allowing
176 model-based RL to scale to complex high dimensional tasks [17, 28]. However, despite the incredible
177 progress, none of these prior works account for the bias introduced by using learned models for
178 planning.

179 Classically, the field of robust control has looked at designing controllers for with uncertain transition
180 models [3, 9]. More relevantly, [1] incorporates bias correction for approximate transition models
181 to learn better policies. However, they only focus on approximation error which arises due to the
182 use of simple linear models. More recently, [19] discusses the mismatch which arises between the
183 training objectives and the control objectives, and how that discrepancy can lead to biases in model
184 learning. In contrast to these works, our work focuses on the sampling bias introduced despite using
185 flexible function approximators like neural networks. To some extent, our work also addresses the
186 bias discussed in [19], as our objective encourages the transition model to fit better to trajectories
187 with higher return.

188 Prior work has considered estimating value functions using transition models [15, 11, 5]. However,
189 our work addresses the biased values generated by the use of learned transition models. To that extent,
190 we rely on classifier-based density ratio estimation techniques [33], which have been employed in
191 generative modelling [36, 13]. [10] use similar principles to correct for differences in dynamics
192 when transferring experience between environments, where the classifier’s corrections are realized
193 by adjusting the reward appropriately before applying model-free RL. To the best of our knowledge,
194 ours is first work which raises and addresses the issue of sampling bias in context of model-based RL.

195 6 Conclusion

196 We introduce a framework for mitigating the impact of model error on planning in model-based RL.
197 Our framework trains discriminators to correct value estimation bias, while also introducing a novel
198 model training objective for minimizing value estimation variance during planning. Empirically,
199 we find that these modifications improve agent performance in environments with stochastic and
200 multi-modal dynamics that pose a challenge for standard model fitting techniques. By improving
201 planning in environments that are difficult to model properly, this framework is a step towards more
202 efficient and performant reinforcement learning in complex, real world problems.

203 References

- 204 [1] P. Abbeel, M. Quigley, and A. Y. Ng. Using inaccurate models in reinforcement learning. In
205 *Proceedings of the 23rd international conference on Machine learning*, pages 1–8, 2006.
- 206 [2] S. Arora, A. Risteski, and Y. Zhang. Do gans learn the distribution? some theory and empirics.
207 In *International Conference on Learning Representations*, 2018.
- 208 [3] J. A. Bagnell, A. Y. Ng, and J. G. Schneider. Solving uncertain markov decision processes.
209 2001.
- 210 [4] C. M. Bishop. Mixture density networks. 1994.
- 211 [5] J. Buckman, D. Hafner, G. Tucker, E. Brevdo, and H. Lee. Sample-efficient reinforcement learn-
212 ing with stochastic ensemble value expansion. In *Advances in Neural Information Processing*
213 *Systems*, pages 8224–8234, 2018.

- 214 [6] K. Chua, R. Calandra, R. McAllister, and S. Levine. Deep reinforcement learning in a handful
215 of trials using probabilistic dynamics models. In *Advances in Neural Information Processing*
216 *Systems*, pages 4754–4765, 2018.
- 217 [7] M. Deisenroth and C. E. Rasmussen. Pilco: A model-based and data-efficient approach to policy
218 search. In *Proceedings of the 28th International Conference on machine learning (ICML-11)*,
219 pages 465–472, 2011.
- 220 [8] M. P. Deisenroth, D. Fox, and C. E. Rasmussen. Gaussian processes for data-efficient learning
221 in robotics and control. *IEEE transactions on pattern analysis and machine intelligence*,
222 37(2):408–423, 2013.
- 223 [9] L. El Ghaoui and A. Nilim. Robust solutions to markov decision problems with uncertain
224 transition matrices. *Operations Research*, 53(5):780–798, 2005.
- 225 [10] B. Eysenbach, S. Asawa, S. Chaudhari, R. Salakhutdinov, and S. Levine. Off-dynamics reinforce-
226 ment learning: Training for transfer with domain classifiers. *arXiv preprint arXiv:2006.13916*,
227 2020.
- 228 [11] V. Feinberg, A. Wan, I. Stoica, M. I. Jordan, J. E. Gonzalez, and S. Levine. Model-based value
229 estimation for efficient model-free reinforcement learning. *arXiv preprint arXiv:1803.00101*,
230 2018.
- 231 [12] Y. Gal, R. McAllister, and C. E. Rasmussen. Improving pilco with bayesian neural network
232 dynamics models. In *Data-Efficient Machine Learning workshop, ICML*, volume 4, page 34,
233 2016.
- 234 [13] A. Grover, J. Song, A. Kapoor, K. Tran, A. Agarwal, E. J. Horvitz, and S. Ermon. Bias correction
235 of learned generative models using likelihood-free importance weighting. In *Advances in Neural*
236 *Information Processing Systems*, pages 11058–11070, 2019.
- 237 [14] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. Learning latent
238 dynamics for planning from pixels. In *International Conference on Machine Learning*, pages
239 2555–2565. PMLR, 2019.
- 240 [15] N. Heess, G. Wayne, D. Silver, T. Lillicrap, T. Erez, and Y. Tassa. Learning continuous control
241 policies by stochastic value gradients. In *Advances in Neural Information Processing Systems*,
242 pages 2944–2952, 2015.
- 243 [16] M. Janner, J. Fu, M. Zhang, and S. Levine. When to trust your model: Model-based policy
244 optimization. In *Advances in Neural Information Processing Systems*, pages 12519–12530,
245 2019.
- 246 [17] L. Kaiser, M. Babaeizadeh, P. Milos, B. Osinski, R. H. Campbell, K. Czechowski, D. Erhan,
247 C. Finn, P. Kozakowski, S. Levine, et al. Model-based reinforcement learning for atari. *arXiv*
248 *preprint arXiv:1903.00374*, 2019.
- 249 [18] V. Kumar, E. Todorov, and S. Levine. Optimal control with learned local models: Application
250 to dexterous manipulation. In *2016 IEEE International Conference on Robotics and Automation*
251 *(ICRA)*, pages 378–383. IEEE, 2016.
- 252 [19] N. Lambert, B. Amos, O. Yadan, and R. Calandra. Objective mismatch in model-based
253 reinforcement learning. *arXiv preprint arXiv:2002.04523*, 2020.
- 254 [20] E. Langlois, S. Zhang, G. Zhang, P. Abbeel, and J. Ba. Benchmarking model-based reinforce-
255 ment learning. *arXiv preprint arXiv:1907.02057*, 2019.
- 256 [21] S. Levine, C. Finn, T. Darrell, and P. Abbeel. End-to-end training of deep visuomotor policies.
257 *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- 258 [22] W. Li and E. Todorov. Iterative linear quadratic regulator design for nonlinear biological
259 movement systems. In *ICINCO (1)*, pages 222–229, 2004.
- 260 [23] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra.
261 Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.

- 262 [24] K. Lowrey, A. Rajeswaran, S. Kakade, E. Todorov, and I. Mordatch. Plan online, learn offline:
 263 Efficient learning and exploration via model-based control. *arXiv preprint arXiv:1811.01848*,
 264 2018.
- 265 [25] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves,
 266 M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep rein-
 267 forcement learning. *nature*, 518(7540):529–533, 2015.
- 268 [26] I. Mordatch, E. Todorov, and Z. Popović. Discovery of complex behaviors through contact-
 269 invariant optimization. *ACM Transactions on Graphics (TOG)*, 31(4):1–8, 2012.
- 270 [27] A. Nagabandi, G. Kahn, R. S. Fearing, and S. Levine. Neural network dynamics for model-
 271 based deep reinforcement learning with model-free fine-tuning. In *2018 IEEE International*
 272 *Conference on Robotics and Automation (ICRA)*, pages 7559–7566. IEEE, 2018.
- 273 [28] A. Nagabandi, K. Konolige, S. Levine, and V. Kumar. Deep dynamics models for learning
 274 dexterous manipulation. In *Conference on Robot Learning*, pages 1101–1112, 2020.
- 275 [29] S. Nair, S. Savarese, and C. Finn. Goal-aware prediction: Learning to model what matters.
 276 *arXiv preprint arXiv:2007.07170*, 2020.
- 277 [30] M. L. Puterman. Markov decision processes. *Handbooks in operations research and manage-*
 278 *ment science*, 2:331–434, 1990.
- 279 [31] R. Y. Rubinstein and D. P. Kroese. *The cross-entropy method: a unified approach to combinato-*
 280 *rial optimization, Monte-Carlo simulation and machine learning*. Springer Science & Business
 281 Media, 2013.
- 282 [32] A. Sharma, S. Gu, S. Levine, V. Kumar, and K. Hausman. Dynamics-aware unsupervised
 283 discovery of skills. *arXiv preprint arXiv:1907.01657*, 2019.
- 284 [33] M. Sugiyama, T. Suzuki, and T. Kanamori. *Density ratio estimation in machine learning*.
 285 Cambridge University Press, 2012.
- 286 [34] R. S. Sutton. Integrated architectures for learning, planning, and reacting based on approximating
 287 dynamic programming. In *Machine learning proceedings 1990*, pages 216–224. Elsevier, 1990.
- 288 [35] Y. Tassa, T. Erez, and E. Todorov. Synthesis and stabilization of complex behaviors through
 289 online trajectory optimization. In *2012 IEEE/RSJ International Conference on Intelligent*
 290 *Robots and Systems*, pages 4906–4913. IEEE, 2012.
- 291 [36] D. Tran, R. Ranganath, and D. Blei. Hierarchical implicit models and likelihood-free variational
 292 inference. In *Advances in Neural Information Processing Systems*, pages 5523–5533, 2017.
- 293 [37] T. Wang and J. Ba. Exploring model-based planning with policy networks. *arXiv preprint*
 294 *arXiv:1906.08649*, 2019.
- 295 [38] G. Williams, A. Aldrich, and E. Theodorou. Model predictive path integral control using
 296 covariance variable importance sampling. *arXiv preprint arXiv:1509.01149*, 2015.