
Accelerating Online Reinforcement Learning with Offline Datasets

Ashvin Nair
UC Berkeley

Murtaza Dalal
UC Berkeley

Abhishek Gupta
UC Berkeley

Sergey Levine
UC Berkeley

Abstract

Reinforcement learning provides an appealing formalism for learning control policies from experience. However, the classic active formulation of reinforcement learning necessitates a lengthy active exploration process for each behavior, making it difficult to apply in real-world settings. If we can instead allow reinforcement learning to effectively use previously collected data to aid the online learning process, where the data could be expert demonstrations or more generally any prior experience, we could make reinforcement learning a substantially more practical tool. While a number of recent methods have sought to learn offline from previously collected data, it remains exceptionally difficult to train a policy with offline data and improve it further with online reinforcement learning. In this paper we systematically analyze why this problem is so challenging, and propose a novel algorithm that combines sample-efficient dynamic programming with maximum likelihood policy updates, providing a simple and effective framework that is able to leverage large amounts of offline data and then quickly perform online fine-tuning of reinforcement learning policies. We show that our method enables rapid learning of skills with a combination of prior demonstration data and online experience across a suite of difficult dexterous manipulation and benchmark tasks.

1 Introduction

Learning models that generalize effectively to complex open-world settings, from image recognition [8] to natural language processing [3], relies on large, high-capacity models and large, diverse, and representative datasets. Leveraging this recipe for reinforcement learning (RL) has the potential to yield powerful policies for real-world control applications such as robotics. However, while deep RL algorithms enable the use of large models, the use of large datasets for real-world RL is conceptually challenging. Most RL algorithms collect new data online every time a new policy is learned, which limits the size and diversity of the datasets for RL. In the same way that powerful models in computer vision and NLP are often pre-trained on large, general-purpose datasets and then fine-tuned on task-specific data, RL policies that generalize effectively to open-world settings will need to be able to incorporate large amounts of prior data effectively into the learning process, while still collecting additional data online for the task at hand.

For data-driven reinforcement learning, offline datasets consist of trajectories of states, actions and associated rewards. This data can potentially come from demonstrations for the desired task [16, 2], suboptimal policies [6], demonstrations for related tasks [21], or even just random exploration in the environment. Depending on the quality of the data that is provided, useful knowledge can be extracted about the dynamics of the world, about the task being solved, or both. Effective data-driven methods for deep reinforcement learning should be able to use this data to pre-train offline while improving with online fine-tuning.

Since this prior data can come from a variety of sources, we require an algorithm that does not utilize different types of data in any privileged way. For example, prior methods that incorporate demonstrations into RL directly aim to mimic these demonstrations [11], which is desirable when the demonstrations are known to be optimal, but can cause undesirable bias when the prior data is not optimal. While prior methods for fully offline RL provide a mechanism for utilizing offline data [5, 9], as we will show in our experiments, such methods generally are not effective for fine-tuning with online data as they are often too conservative. In effect, prior methods require us to choose: Do we assume prior data is optimal or not? Do we use only offline data, or only online data? To make it feasible to learn policies for open-world settings, we need algorithms that contain all of the aforementioned qualities.

In this work, we study how to build RL algorithms that are effective for pre-training from a variety of off-policy datasets, but also well suited to continuous improvement with online data collection. We systematically analyze the challenges with using standard off-policy RL algorithms [7, 9, 1] for this problem, and introduce a simple actor critic algorithm that elegantly bridges data-driven pre-training from offline data and improvement with online data collection. Our method, which uses dynamic programming to train a critic but a supervised update to train a constrained actor, combines the best of supervised learning and actor-critic algorithms. Dynamic programming can leverage off-policy data and enable sample-efficient learning. The simple supervised actor update implicitly enforces a constraint that mitigates the effects of out-of-distribution actions when learning from offline data [5, 9], while avoiding overly conservative updates. We evaluate our algorithm on a wide variety of robotic control and benchmark tasks across three simulated domains: dexterous manipulation, tabletop manipulation, and MuJoCo control tasks. We see that our algorithm, Advantage Weighted Actor Critic (AWAC), is able to quickly learn successful policies on difficult tasks with high action dimension and binary sparse rewards, significantly better than prior methods for off-policy and offline reinforcement learning. Moreover, we see that AWAC can utilize different types of prior data: demonstrations, suboptimal data, and random exploration data.

2 Advantage Weighted Actor Critic: A Simple Algorithm for Fine-tuning from Offline Datasets

In this section, we will describe the advantage weighted actor-critic (AWAC) algorithm, which trains an off-policy critic and an actor with an *implicit* policy constraint. AWAC follows the standard paradigm for actor-critic algorithms, with a policy evaluation step to learn Q^π and a policy improvement step to update π . AWAC uses off-policy temporal-difference learning to estimate Q^π in the policy evaluation step, and a unique policy improvement update that is able to obtain the benefits of offline RL algorithms at training from prior datasets, while avoiding the overly conservative behavior of offline algorithms. We describe the policy improvement step in AWAC below, and summarize the entire algorithm thereafter.

Policy improvement for AWAC proceeds by learning a policy that maximizes the value of the critic learned in the policy evaluation step via TD bootstrapping. At iteration k , AWAC therefore optimizes the policy to maximize the estimated Q-function $Q^{\pi_k}(\mathbf{s}, \mathbf{a})$ at every state, while constraining it to stay close to the actions observed in the data, similar to prior offline RL methods, though this constraint will be enforced differently. Using advantages instead of Q-values, we can write this optimization as:

$$\pi_{k+1} = \arg \max_{\pi \in \Pi} \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\mathbf{s})} [A^{\pi_k}(\mathbf{s}, \mathbf{a})] \text{ s.t. } D_{\text{KL}}(\pi(\cdot|\mathbf{s}) || \pi_\beta(\cdot|\mathbf{s})) \leq \epsilon. \quad (1)$$

We first derive the solution to the constrained optimization in Equation 1 to obtain a non-parametric closed form for the actor. This solution is then projected onto the parametric policy class *without* any explicit behavior model. The analytic solution to Equation 1 can be obtained by enforcing the KKT conditions [13, 14, 12]. The Lagrangian is:

$$\mathcal{L}(\pi, \lambda) = \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\mathbf{s})} [A^{\pi_k}(\mathbf{s}, \mathbf{a})] + \lambda(\epsilon - D_{\text{KL}}(\pi(\cdot|\mathbf{s}) || \pi_\beta(\cdot|\mathbf{s}))), \quad (2)$$

and the closed form solution to this problem is $\pi^*(\mathbf{a}|\mathbf{s}) \propto \frac{1}{Z(\mathbf{s})} \pi_\beta(\mathbf{a}|\mathbf{s}) \exp\left(\frac{1}{\lambda} A^{\pi_k}(\mathbf{s}, \mathbf{a})\right)$. When using function approximators, such as deep neural networks as we do in our implementation, we need to project the non-parametric solution into our policy space. For a policy π_θ with parameters θ , this can be done by minimizing the KL divergence of π_θ from the optimal non-parametric solution π^*

under the data distribution $\rho_{\pi_\beta}(\mathbf{s})$:

$$\arg \min_{\theta} \mathbb{E}_{\rho_{\pi_\beta}(\mathbf{s})} [D_{\text{KL}}(\pi^*(\cdot|\mathbf{s})||\pi_\theta(\cdot|\mathbf{s}))] = \arg \min_{\theta} \mathbb{E}_{\rho_{\pi_\beta}(\mathbf{s})} \left[\mathbb{E}_{\pi^*(\cdot|\mathbf{s})} [-\log \pi_\theta(\cdot|\mathbf{s})] \right] \quad (3)$$

Note that the parametric policy could be projected with either direction of KL divergence. Choosing the reverse KL results in explicit penalty methods [20] that rely on evaluating the density of a learned behavior model. Instead, by using forward KL, we can sample directly from β :

$$\theta_{k+1} = \arg \max_{\theta} \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \beta} \left[\log \pi_\theta(\mathbf{a}|\mathbf{s}) \exp \left(\frac{1}{\lambda} A^{\pi_k}(\mathbf{s}, \mathbf{a}) \right) \right]. \quad (4)$$

This actor update amounts to weighted maximum likelihood (i.e., supervised learning), where the targets are obtained by re-weighting the state-action pairs observed in the current dataset by the predicted advantages from the learned critic, *without* explicitly learning any parametric behavior model, simply sampling (s, a) from the replay buffer β . See Appendix B.2 for a more detailed derivation and Appendix B.3 for specific implementation details.

Avoiding explicit behavior modeling. Note that the update in Equation 4 completely avoids any modeling of the previously observed data β with a parametric model. By avoiding any explicit learning of the behavior model AWAC is far less conservative than methods which fit a model $\hat{\pi}_\beta$ explicitly, and better incorporates new data during online fine-tuning, as seen from our results in Section 3. This derivation is related to AWR [12], with the main difference that AWAC uses an off-policy Q-function Q^π to estimate the advantage, which greatly improves efficiency and even final performance (see results in Section 6.1). The update also resembles ABM-MPO, but ABM-MPO *does* require modeling the behavior policy which can lead to poor fine-tuning. In Section 6.1, AWAC outperforms ABM-MPO on a range of challenging tasks.

Policy evaluation. During policy evaluation, we estimate the action-value $Q^\pi(\mathbf{s}, \mathbf{a})$ for the current policy π . We utilize a standard temporal difference learning scheme for policy evaluation [7, 4], by minimizing the Bellman error. This enables us to learn very efficiently from off-policy data. This is particularly important in our problem setting to effectively use the offline dataset, and allows us to significantly outperform alternatives using Monte-Carlo evaluation or TD(λ) to estimate returns [12].

Algorithm summary. The full AWAC algorithm for offline RL with online fine-tuning is summarized in Algorithm 1. In a practical implementation, we can parameterize the actor and the critic by neural networks and perform SGD updates, alternating between Eqn. 4 and Bellman updates. Specific details are provided in Appendix B.3. As we will show in our experiments, the specific design choices described above enable AWAC to excel fine-tuning after pretraining. AWAC ensures data efficiency with off-policy critic estimation via bootstrapping, and avoids offline bootstrap error with a constrained actor update. By avoiding explicit modeling of the behavior policy, AWAC avoids overly conservative updates.

Algorithm 1 Advantage Weighted AC

```

1: Dataset  $\mathcal{D} = \{(\mathbf{s}, \mathbf{a}, \mathbf{s}', r)_j\}$ 
2: Initialize buffer  $\beta = \mathcal{D}$ 
3: Initialize  $\pi_\theta, Q_\phi$ 
4: for iteration  $i = 1, 2, \dots$  do
5:   Sample batch  $(\mathbf{s}, \mathbf{a}, \mathbf{s}', r) \sim \beta$ 
6:   Update  $\phi$  with Bellman eqn.
7:   Update  $\theta$  according to Eqn. 4
8:   if  $i > \text{num\_offline\_steps}$  then
9:      $\tau_1, \dots, \tau_K \sim p_{\pi_\theta}(\tau)$ 
10:     $\beta \leftarrow \beta \cup \{\tau_1, \dots, \tau_K\}$ 
11:   end if
12: end for
```

3 Experimental Evaluation

In our experiments, we first compare our method against prior methods in the offline training and fine-tuning setting. We show that we can learn difficult, high-dimensional, sparse reward dexterous manipulation problems from human demonstrations and off-policy data. We then evaluate our method with suboptimal prior data generated by a random controller. Finally, we study why prior methods struggle in this setting by analyzing their performance on benchmark MuJoCo tasks, and conduct further experiments to understand where the difficulty lies. Videos and further experimental details can also be found at sites.google.com/view/awac-anonymous

6.1) Comparative Evaluation on Dexterous Manipulation Tasks. We aim to study tasks representative of the difficulties of real-world robot learning, where offline learning and online fine-tuning are most relevant. One such setting is the suite of dexterous manipulation tasks proposed by Rajeswaran et al. [15]. These tasks involve complex manipulation skills using a 28-DoF five-fingered hand in the MuJoCo simulator [18] shown in Figure 1: in-hand rotation of a pen, opening a door by unlatching

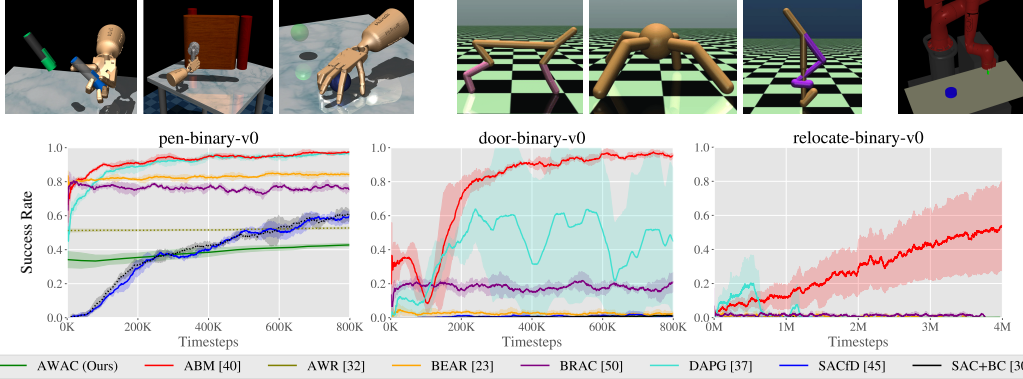


Figure 1: Comparative evaluation on the dexterous manipulation tasks. These tasks are difficult due to their high action dimensionality and reward sparsity. We see that AWAC is able to learn these tasks with little online data collection required (100K samples \approx 16 minutes of equivalent real-world interaction time). Meanwhile, most prior methods are not able to solve the harder two tasks: door opening and object relocation.

the handle, and picking up a sphere and relocating it to a target location. These environments exhibit many challenges: high dimensional action spaces, complex manipulation physics with many intermittent contacts, and randomized hand and object positions. The reward functions in these environments are binary 0-1 rewards for task completion.¹ Rajeswaran et al. [15] provide 25 human demonstrations for each task, which are not fully optimal but do solve the task. Since this dataset is very small, we generated another 500 trajectories with behavior cloning.

First, we compare our method on the dexterous manipulation tasks described earlier against prior methods for off-policy learning, offline learning, and bootstrapping from demonstrations. Specific implementation details are discussed in Appendix B.4. The results are shown in Fig. 1. Our method is able to leverage the prior data to quickly attain good performance, and the efficient off-policy actor-critic component of our approach fine-tunes much more quickly than demonstration augmented policy gradient (DAPG), the method proposed by Rajeswaran et al. [15]. For example, our method solves the pen task in 120K timesteps, the equivalent of just 20 minutes of online interaction. While the baseline comparisons and ablations are able to make some amount of progress on the pen task, alternative off-policy RL and offline RL algorithms are largely unable to solve the door and relocate task in the time-frame considered. We find that the design decisions to use off-policy critic estimation allow AWAC to significantly outperform AWR [12] while the implicit behavior modeling allows AWAC to significantly outperform ABM [17], although ABM does make some progress. Rajeswaran et al. [15] show that DAPG can eventually solve these tasks with more reward information, but this highlights the weakness of on-policy methods in sparse reward scenarios. Similar results hold true in benchmark Gym environments; due to space constraints, we present these results in the appendix.

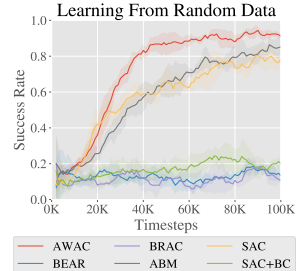


Figure 2: Comparison of fine-tuning from an initial dataset of suboptimal data on a robot pushing task.

6.2) Fine-Tuning from Random Policy Data. An advantage of using off-policy RL for reinforcement learning is that we can also incorporate suboptimal data, rather than only demonstrations. In this experiment, we evaluate on a simulated tabletop pushing environment with a Sawyer robot (shown in Fig 1), described further in Appendix B.1. To study the potential to learn from suboptimal data, we use an off-policy dataset of 500 trajectories generated by a random process. The task is to push an object to a target location in a 40cm x 20cm goal space.

The results are shown in Figure 2. We see that while many methods begin at the same initial performance, AWAC learns the fastest online and is actually able to make use of the offline dataset effectively as opposed to some methods which are completely unable to learn.

¹Rajeswaran et al. [15] use a combination of task completion factors as the sparse reward. For instance, in the door task, the sparse reward as a function of the door position d was $r = 10\mathbb{1}_{d>1.35} + 8\mathbb{1}_{d>1.0} + 2\mathbb{1}_{d>1.2} - 0.1\|d - 1.57\|_2$. We only use the success measure $r = \mathbb{1}_{d>1.4}$, which is substantially more difficult.

References

- [1] Abdolmaleki, A., Springenberg, J. T., Tassa, Y., Munos, R., Heess, N., and Riedmiller, M. Maximum a Posteriori Policy Optimisation. In *International Conference on Learning Representations (ICLR)*, pp. 1–19, 2018.
- [2] Atkeson, C. G. and Schaal, S. Robot Learning From Demonstration. In *International Conference on Machine Learning (ICML)*, 1997.
- [3] Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Association for Computational Linguistics (ACL)*, oct 2019.
- [4] Fujimoto, S., van Hoof, H., and Meger, D. Addressing Function Approximation Error in Actor-Critic Methods. *International Conference on Machine Learning (ICML)*, 2018.
- [5] Fujimoto, S., Meger, D., and Precup, D. Off-Policy Deep Reinforcement Learning without Exploration. In *International Conference on Machine Learning (ICML)*, dec 2019.
- [6] Gao, Y., Xu, H., Lin, J., Yu, F., Levine, S., and Darrell, T. Reinforcement learning from imperfect demonstrations. *CoRR*, abs/1802.05313, 2018.
- [7] Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. In *International Conference on Machine Learning*, 2018.
- [8] Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems (NIPS)*, pp. 1097–1105, 2012.
- [9] Kumar, A., Fu, J., Tucker, G., and Levine, S. Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction. In *Neural Information Processing Systems (NeurIPS)*, jun 2019.
- [10] Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. Continuous control with deep reinforcement learning. In *International Conference on Learning Representations (ICLR)*, 2016. ISBN 0-7803-3213-X. doi: 10.1613/jair.301.
- [11] Nair, A., McGrew, B., Andrychowicz, M., Zaremba, W., and Abbeel, P. Overcoming Exploration in Reinforcement Learning with Demonstrations. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2018.
- [12] Peng, X. B., Kumar, A., Zhang, G., and Levine, S. Advantage-Weighted Regression: Simple and Scalable Off-Policy Reinforcement Learning. sep 2019.
- [13] Peters, J. and Schaal, S. Reinforcement Learning by Reward-weighted Regression for Operational Space Control. In *International Conference on Machine Learning*, 2007.
- [14] Peters, J., Mülling, K., and Altün, Y. Relative Entropy Policy Search. In *AAAI Conference on Artificial Intelligence*, pp. 1607–1612, 2010.
- [15] Rajeswaran, A., Kumar, V., Gupta, A., Schulman, J., Todorov, E., and Levine, S. Learning Complex Dexterous Manipulation with Deep Reinforcement Learning and Demonstrations. In *Robotics: Science and Systems*, 2018.
- [16] Schaal, S. Learning from demonstration. In *Advances in Neural Information Processing Systems (NeurIPS)*, number 9, pp. 1040–1046, 1997. ISBN 1558604863. doi: 10.1016/j.robot.2004.03.001.
- [17] Siegel, N. Y., Springenberg, J. T., Berkenkamp, F., Abdolmaleki, A., Neunert, M., Lampe, T., Hafner, R., Heess, N., and Riedmiller, M. Keep doing what worked: Behavioral modelling priors for offline reinforcement learning, 2020.
- [18] Todorov, E., Erez, T., and Tassa, Y. MuJoCo: A physics engine for model-based control. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 5026–5033, 2012. ISBN 9781467317375. doi: 10.1109/IROS.2012.6386109.
- [19] Večerík, M., Hester, T., Scholz, J., Wang, F., Pietquin, O., Piot, B., Heess, N., Rothörl, T., Lampe, T., and Riedmiller, M. Leveraging Demonstrations for Deep Reinforcement Learning on Robotics Problems with Sparse Rewards. *CoRR*, abs/1707.0, 2017.
- [20] Wu, Y., Tucker, G., and Nachum, O. Behavior Regularized Offline Reinforcement Learning. In *International Conference on Learning Representations (ICLR)*, nov 2020.
- [21] Zhou, A., Jang, E., Kappler, D., Herzog, A., Khansari, M., Wohlhart, P., Bai, Y., Kalakrishnan, M., Levine, S., and Finn, C. Watch, try, learn: Meta-learning from demonstrations and reward. *CoRR*, abs/1906.03352, 2019.

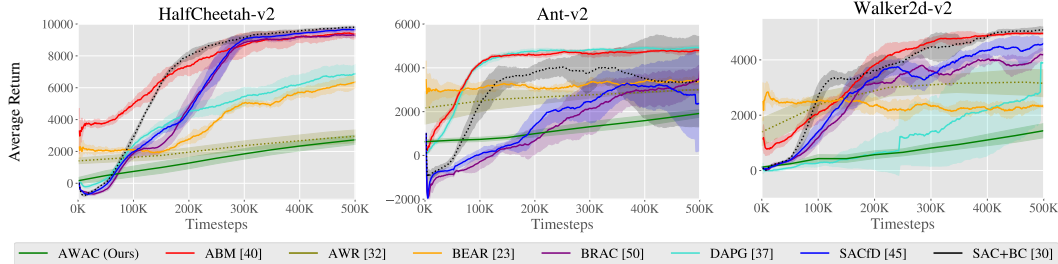


Figure 3: Comparison of our method and prior methods on standard MuJoCo benchmark tasks.

A Appendix

B analysis on MuJoCo Benchmarks from Prior Data.

Since the dexterous manipulation environments are challenging to solve, we provide a comparative evaluation on MuJoCo benchmark tasks for analysis. On these simpler problems, many prior methods are able to learn, but it allows us to understand more precisely which design decisions are crucial. For each task, we collect 15 demonstration trajectories using a pre-trained expert on each task, and 100 trajectories of off-policy data by rolling out a behavioral cloned policy trained on the demonstrations. The same data is made available to all methods. The results are presented in Figure 3. AWAC is consistently the best-performing method, but several other methods show reasonable performance.

Data efficiency. The two methods that do not estimate Q^π are DAPG [1] and AWR [12]. Across all three tasks, we see that these methods are somewhat worse offline than the best performing offline methods, and exhibit steady but slow improvement.

Bootstrap error in offline off-policy learning. For SAC [7], across all three tasks, we see that the offline performance at epoch 0 is generally poor. Due to the data in the replay buffer, SAC with prior data does learn faster than from scratch, but AWAC is faster to solve the tasks in general. SAC with additional data in the replay buffer is similar to the approach proposed by Večerík et al. [19]. SAC+BC reproduces Nair et al. [11] but uses SAC instead of DDPG [10] as the underlying RL algorithm. We find that these algorithms exhibit a characteristic dip at the start of learning.

Conservative online learning. Finally, we consider conservative offline algorithms: ABM [17], BEAR [9], and BRAC [20]. We found that BRAC performs similarly to SAC for working hyperparameters. BEAR trains well offline - on Ant and Walker2d, BEAR significantly outperforms prior methods before online experience. However, online improvement is slow for BEAR and the final performance across all three tasks is much lower than AWAC. The closest in performance to our method is ABM, which is comparable on Ant-v2, but much slower on other domains.

B.1 Environment-Specific Details

We evaluate our method on three domains: dexterous manipulation environments, Sawyer manipulation environments, and MuJoCo benchmark environments. In the following sections we describe specific details.

B.1.1 Dexterous Manipulation Environments

These environments are modified from those proposed by by Rajeswaran et al. [15], and available in this repository.

pen-binary-v0. The task is to spin a pen into a given orientation. The action dimension is 24 and the observation dimension is 45. Let the position and orientation of the pen be denoted by x_p and x_o respectively, and the desired position and orientation be denoted by d_p and d_o respectively. The reward function is $r = \mathbb{1}_{|x_p - d_p| \leq 0.075} \mathbb{1}_{|x_o - d_o| \leq 0.95} - 1$. In Rajeswaran et al. [15], the episode was terminated when the pen fell out of the hand; we did not include this early termination condition.

door-binary-v0. The task is to open a door, which requires first twisting a latch. The action dimension is 28 and the observation dimension is 39. Let d denote the angle of the door. The reward function is $r = \mathbb{1}_{d>1.4} - 1$.

relocate-binary-v0. The task is to relocate an object to a goal location. The action dimension is 30 and the observation dimension is 39. Let x_p denote the object position and d_p denote the desired position. The reward is $r = \mathbb{1}_{|x_p - d_p| \leq 0.1} - 1$.

B.1.2 Sawyer Manipulation Environment

SawyerPush-v0. This environment is included in the Multiworld library. The task is to push a puck to a goal position in a 40cm x 20cm, and the reward function is the negative distance between the puck and goal position. When using this environment, we use hindsight experience replay for goal-conditioned reinforcement learning. The random dataset for prior data was collected by rolling out an Ornstein-Uhlenbeck process with $\theta = 0.15$ and $\sigma = 0.3$.

B.2 Algorithm Derivation Details

The full optimization problem we solve, given the previous off-policy advantage estimate A^{π_k} and buffer distribution π_β is:

$$\pi_{k+1} = \arg \max_{\pi \in \Pi} \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\mathbf{s})} [A^{\pi_k}(\mathbf{s}, \mathbf{a})] \quad (5)$$

$$\text{s.t. } D_{\text{KL}}(\pi(\cdot|\mathbf{s}) || \pi_\beta(\cdot|\mathbf{s})) \leq \epsilon \quad (6)$$

$$\int_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{s}) d\mathbf{a} = 1. \quad (7)$$

Our derivation follows Peters et al. [14] and Peng et al. [12]. The analytic solution for the constrained optimization problem above can be obtained by enforcing the KKT conditions. The Lagrangian is:

$$\mathcal{L}(\pi, \lambda, \alpha) = \mathbb{E}_{\mathbf{a} \sim \pi(\cdot|\mathbf{s})} [A^{\pi_k}(\mathbf{s}, \mathbf{a})] + \lambda(\epsilon - D_{\text{KL}}(\pi(\cdot|\mathbf{s}) || \pi_\beta(\cdot|\mathbf{s}))) + \alpha(1 - \int_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{s}) d\mathbf{a}). \quad (8)$$

Differentiating with respect to π gives:

$$\frac{\partial \mathcal{L}}{\partial \pi} = A^{\pi_k}(\mathbf{s}, \mathbf{a}) - \lambda \log \pi_\beta(\mathbf{a}|\mathbf{s}) + \lambda \log \pi(\mathbf{a}|\mathbf{s}) + \lambda - \alpha. \quad (9)$$

Setting $\frac{\partial \mathcal{L}}{\partial \pi}$ to zero and solving for π gives the closed form solution to this problem:

$$\pi^*(\mathbf{a}|\mathbf{s}) = \frac{1}{Z(\mathbf{s})} \pi_\beta(\mathbf{a}|\mathbf{s}) \exp \left(\frac{1}{\lambda} A^{\pi_k}(\mathbf{s}, \mathbf{a}) \right), \quad (10)$$

Next, we project the solution into the space of parametric policies. For a policy π_θ with parameters θ , this can be done by minimizing the KL divergence of π_θ from the optimal non-parametric solution π^* under the data distribution $\rho_{\pi_\beta}(\mathbf{s})$:

$$\arg \min_{\theta} \mathbb{E}_{\rho_{\pi_\beta}(\mathbf{s})} [D_{\text{KL}}(\pi^*(\cdot|\mathbf{s}) || \pi_\theta(\cdot|\mathbf{s}))] = \arg \min_{\theta} \mathbb{E}_{\rho_{\pi_\beta}(\mathbf{s})} \left[\mathbb{E}_{\pi^*(\cdot|\mathbf{s})} [-\log \pi_\theta(\cdot|\mathbf{s})] \right] \quad (11)$$

Note that in the projection step, the parametric policy could be projected with either direction of KL divergence. However, choosing the reverse KL direction has a key advantage: it allows us to optimize θ as a maximum likelihood problem with an expectation over data $s, a \sim \beta$, rather than sampling actions from the policy that may be out of distribution for the Q function. In our experiments we show that this decision is vital for stable off-policy learning.

Furthermore, assume discrete policies with a minimum probably density of $\pi_\theta \geq \alpha_\theta$. Then the upper bound:

$$D_{\text{KL}}(\pi^* || \pi_\theta) \leq \frac{2}{\alpha_\theta} D_{\text{TV}}(\pi^*, \pi_\theta)^2 \quad (12)$$

$$\leq \frac{1}{\alpha_\theta} D_{\text{KL}}(\pi_\theta || \pi^*) \quad (13)$$

holds by the Pinsker's inequality, where D_{TV} denotes the total variation distance between distributions. Thus minimizing the reverse KL also bounds the forward KL. Note that we can control the minimum α if desired by applying Laplace smoothing to the policy.

Hyper-parameter	Value
Training Batches Per Timestep	1
Exploration Noise	None (stochastic policy)
RL Batch Size	1024
Discount Factor	0.99
Reward Scaling	1
Replay Buffer Size	1000000
Number of pretraining steps	25000
Policy Hidden Sizes	[256, 256, 256, 256]
Policy Hidden Activation	ReLU
Policy Weight Decay	10^{-4}
Policy Learning Rate	3×10^{-4}
Q Hidden Sizes	[256, 256, 256, 256]
Q Hidden Activation	ReLU
Q Weight Decay	0
Q Learning Rate	3×10^{-4}
Target Network τ	5×10^{-3}

Table 1: Hyper-parameters used for RL experiments.

Name	\hat{Q}	Policy Objective	$\hat{\pi}_\beta$?	Constraint
SAC	Q^π	$D_{\text{KL}}(\pi_\theta \bar{Q})$	No	None
SAC + BC	Q^π	Mixed	No	None
BCQ	Q^π	$D_{\text{KL}}(\pi_\theta \bar{Q})$	Yes	Support (ℓ^∞)
BEAR	Q^π	$D_{\text{KL}}(\pi_\theta \bar{Q})$	Yes	Support (MMD)
AWR	Q^β	$D_{\text{KL}}(\bar{Q} \pi_\theta)$	No	Implicit
MPO	Q^π	$D_{\text{KL}}(\bar{Q} \pi_\theta)$	Yes*	Prior
ABM-MPO	Q^π	$D_{\text{KL}}(\bar{Q} \pi_\theta)$	Yes	Learned Prior
DAPG	-	$J(\pi_\theta)$	No	None
BRAC	Q^π	$D_{\text{KL}}(\pi_\theta \bar{Q})$	Yes	Explicit KL penalty
AWAC (Ours)	Q^π	$D_{\text{KL}}(\bar{Q} \pi_\theta)$	No	Implicit

Figure 4: Comparison of prior algorithms that can incorporate prior datasets. See section B.4 for specific implementation details.

B.3 Implementation Details

We implement the algorithm building on top of twin delayed deep deterministic policy gradient (TD3) from [4]. The base hyperparameters are given in table 1.

The policy update is replaced with:

$$\theta_{k+1} = \arg \max_{\theta} \mathbb{E}_{\mathbf{s}, \mathbf{a} \sim \beta} \left[\log \pi_{\theta}(\mathbf{a} | \mathbf{s}) \frac{1}{Z(\mathbf{s})} \exp \left(\frac{1}{\lambda} A^{\pi_k}(\mathbf{s}, \mathbf{a}) \right) \right]. \quad (14)$$

We found that explicitly computing $Z(s) = \int_{\mathbf{a}} \pi_{\beta}(\mathbf{a} | \mathbf{s}) \exp(\frac{1}{\lambda} A^{\pi_k}(\mathbf{s}, \mathbf{a})) d\mathbf{a}$ results in worse performance, so we ignore the effect of $Z(s)$ and empirically find that this results in strong performance both offline and online.

The Lagrange multiplier λ is a hyperparameter. In this work we use $\lambda = 0.3$ for the manipulation environments and $\lambda = 1.0$ for the MuJoCo benchmark environments. One could adaptively learn λ with a dual gradient descent procedure, but this would require access to π_{β} .

As rewards for the dextrous manipulation environments are non-positive, we clamp the Q value for these experiments to be at most zero. We find this stabilizes training slightly.

B.4 Baseline Implementation Details

We used public implementations of prior methods (DAPG, AWR) when available. We implemented the remaining algorithms in our framework, which also allows us to understand the effects of changing individual components of the method. In the section, we describe the implementation details. The full overview of algorithms is given in Figure 4.

Behavior Cloning (BC). This method learns a policy with supervised learning on demonstration data.

Soft Actor Critic (SAC). Using the soft actor critic algorithm from [7], we follow the exact same procedure as our method in order to incorporate prior data, initializing the policy with behavior cloning on demonstrations and adding all prior data to the replay buffer.

Behavior Regularized Actor Critic (BRAC). We implement BRAC as described in [20] by adding policy regularization $\log(\pi_\beta(a|s))$ where π_β is a behavior policy trained with supervised learning on the replay buffer. We add all prior data to the replay buffer before online training.

Advantage Weighted Regression (AWR). Using the advantage weighted regression algorithm from [12], we add all prior data to the replay buffer before online training. We use the implementation provided by Peng et al. [12], with the key difference from our method being that AWR uses TD(λ) on the replay buffer for policy evaluation.

Maximum a Posteriori Policy Optimization (MPO). We evaluate the MPO algorithm presented by Abdolmaleki et al. [1]. Due to a public implementation being unavailable, we modify our algorithm to be as close to MPO as possible. In particular, we change the policy update in Advantage Weighted Actor Critic to be:

$$\theta_i \leftarrow \arg \max_{\theta_i} \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi(a|s)} \left[\log \pi_{\theta_i}(a|s) \exp\left(\frac{1}{\beta} Q^{\pi_\beta}(s, a)\right) \right]. \quad (15)$$

Note that in MPO, actions for the update are sampled from the policy and the Q-function is used instead of advantage for weights. We failed to see offline or online improvement with this implementation in most environments, so we omit this comparison in favor of ABM.

Advantage-Weighted Behavior Model (ABM). We evaluate ABM, the method developed in Siegel et al. [17]. As with MPO, we modify our method to implement ABM, as there is no public implementation of the method. ABM first trains an advantage model $\pi_{\theta_{\text{abm}}}(a|s)$:

$$\theta_{\text{abm}} = \arg \max_{\theta_i} \mathbb{E}_{\tau \sim \mathcal{D}} \left[\sum_{t=1}^{|\tau|} \log \pi_{\theta_{\text{abm}}}(a_t|s_t) f(R(\tau_{t:N}) - \hat{V}(s)) \right]. \quad (16)$$

where f is an increasing non-negative function, chosen to be $f = 1_+$. In place of an advantage computed by empirical returns $R(\tau_{t:N}) - \hat{V}(s)$ we use the advantage estimate computed per transition by the Q value $Q(s, a) - V(s)$. This is favorable for running ABM online, as computing $R(\tau_{t:N}) - \hat{V}(s)$ is similar to AWR, which shows slow online improvement. We then use the policy update:

$$\theta_i \leftarrow \arg \max_{\theta_i} \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi_{\text{abm}}(a|s)} \left[\log \pi_{\theta_i}(a|s) \exp \left(\frac{1}{\lambda} (Q^{\pi_i}(s, a) - V^{\pi_i}(s)) \right) \right]. \quad (17)$$

Additionally, for this method, actions for the update are sampled from a behavior policy trained to match the replay buffer and the value function is computed as $V^\pi(s) = Q^\pi(s, a)$ s.t. $a \sim \pi$.

Demonstration Augmented Policy Gradient (DAPG). We directly utilize the code provided in [15] to compare against our method. Since DAPG is an on-policy method, we only provide the demonstration data to the DAPG code to bootstrap the initial policy from.

Bootstrapping Error Accumulation Reduction (BEAR). We utilize the implementation of BEAR provided in rlkit. We provide the demonstration and off-policy data to the method together. Since the original method only involved training offline, we modify the algorithm to include an online training phase. In general we found that the MMD constraint in the method was too conservative. As a result, in order to obtain the results displayed in our paper, we swept the MMD threshold value and chose the one with the best final performance after offline training with offline fine-tuning.