

622 **A Metrics and Quasimetrics**

623 A metric space (\mathcal{M}, d) is composed of a set \mathcal{M} and a metric $d : \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}^+ \cup \{\infty\}$ that compares
 624 two points in that set. Here \mathbb{R}^+ is the set of non-negative real numbers.

625 **Definition 2.** A metric $d : \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}^+ \cup \{\infty\}$ compares two points in set \mathcal{M} and satisfies the
 626 following axioms $\forall m_1, m_2, m_3 \in \mathcal{M}$:

- 627 • $d(m_1, m_2) = 0 \iff m_1 = m_2$ (identity of indiscernibles)
- 628 • $d(m_1, m_2) = d(m_2, m_1)$ (symmetry)
- 629 • $d(m_1, m_2) \leq d(m_1, m_3) + d(m_3, m_2)$ (triangle inequality)

630 A variation on metrics that is important to this paper is *quasimetrics*.

631 **Definition 3.** A quasimetric [61] is a function that satisfies all the properties of a metric, with the
 632 exception of symmetry $d(m_1, m_2) \neq d(m_2, m_1)$.

633 As an example, consider an MDP where the actions and transition dynamics allow an agent to navigate
 634 from any state to any other state. Let $T(s_2|\pi, s_1)$ be the random variable for the first time-step that
 635 state s_2 is encountered by the agent after starting in state s_1 and following policy π . The time-step
 636 metric d_T^π for this MDP can then be defined as

$$d_T^\pi(s_1, s_2) := \mathbb{E} [T(s_2|\pi, s_1)]$$

637 d_T^π is a quasimetric, since the action space and transition function need not be symmetric, meaning the
 638 expected minimum time needed to go from s_1 to s_2 need not be the same as the expected minimum
 639 time needed to go from s_2 to s_1 . The diameter of an MDP [36, 39] is generally calculated by taking
 640 the maximum time-step distance between over all pairs of states in the MDP either under a random
 641 policy or a policy that travels from any state to any other state in as few steps as possible.

642 **B Optimal Transport and Wasserstein-1 Distance**

643 The theory of optimal transport [69, 13] considers the question of how much work must be done to
 644 transport one distribution to another optimally. More concretely, suppose we have a metric space
 645 (\mathcal{M}, d) where \mathcal{M} is a set and d is a metric on \mathcal{M} . See the definitions of metrics and quasimetrics
 646 in Appendix A. For two distributions μ and ν with finite moments on the set \mathcal{M} , the Wasserstein- p
 647 distance is denoted by:

$$W_p(\mu, \nu) := \inf_{\zeta \in Z(\mu, \nu)} \mathbb{E}_{(X, Y) \sim \zeta} [d(X, Y)^p]^{1/p} \tag{10}$$

648 where Z is the space of all possible couplings between μ and ν . Put another way, Z is the space of
 649 all possible distributions $\zeta \in \Delta(\mathcal{M} \times \mathcal{M})$ whose marginals are μ and ν respectively. Finding this
 650 optimal coupling tells us what is the least amount of work, as measured by d , that needs to be done to
 651 convert μ to ν . This Wasserstein- p distance can then be used as a cost function (negative reward) by
 652 an RL agent to match a given target distribution [70, 17].

653 Finding the ideal coupling (meaning finding the optimal transport plan from one distribution to the
 654 other) which gives us an accurate distance is generally considered intractable. However, if what we
 655 need is an accurate estimate of the Wasserstein distance and not the optimal transport plan (as is the
 656 case when we mean to use this distance as part of our intrinsic reward) we can turn our attention
 657 to the dual form of this distance. The Kantorovich-Rubinstein duality [69] for the Wasserstein-1
 658 distance on a ground metric d is of particular interest and gives us the following equality:

$$W_1(\mu, \nu) = \sup_{\text{Lip}(f) \leq 1} \mathbb{E}_{y \sim \nu} [f(y)] - \mathbb{E}_{x \sim \mu} [f(x)] \tag{11}$$

659 where the supremum is over all 1-Lipschitz functions $f : \mathcal{M} \rightarrow \mathbb{R}$ in the metric space, and the
 660 Lipschitz constant of a function f is defined as:

$$\text{Lip}(f) := \sup \left\{ \frac{|f(y) - f(x)|}{d(x, y)} \mid \forall (x, y) \in \mathcal{M}^2, x \neq y \right\} \quad (12)$$

661 That is, the Lipschitz condition of this function f (called the Kantorovich potential function) is
 662 measured according to the metric d . Recently, Jevtić [37] has shown that this dual formulation where
 663 the constraint on the potential function is a smoothness constraint extends to quasimetric spaces as
 664 well. If defined over a quasimetric space, the Wasserstein distance also has properties of a quasimetric
 665 (specifically, the distances are not necessarily symmetric).

666 If the given metric space is a Euclidean space ($d(x, y) = \|y - x\|_2$), the Lipschitz bound in Equation
 667 2 can be computed locally as a uniform bound on the gradient of f .

$$W_1(\mu, \nu) = \sup_{\|\nabla f\| \leq 1} \mathbb{E}_{y \sim \nu} [f(y)] - \mathbb{E}_{x \sim \mu} [f(x)] \quad (13)$$

668 meaning that f is the solution to an optimization objective with the restriction that $\|\nabla f(x)\| \leq 1$ for
 669 all $x \in \mathcal{M}$. This strong bound on the dual in Euclidean space is the one that has been used most in
 670 recent implementations of the Wasserstein generative adversarial network [3, 29] to regularize the
 671 learning of the discriminator function. Such regularization has been found to be effective for stability
 672 in other adversarial learning approaches such as adversarial imitation learning [25].

673 Practically, the Kantorovich potential function f can be approximated using samples from the two
 674 distributions μ and ν , regularization of the potential function to ensure smoothness, and an expressive
 675 function approximator such as a neural network. A more in depth treatment of the Kantorovich
 676 relaxation and the Kantorovich-Rubinstein duality, as well as their application in metric and Euclidean
 677 spaces using the Wasserstein-1 distance we lay out above, is provided by Peyré and Cuturi [53].

678 Now consider the problem of goal-conditioned reinforcement learning. Here the target distribution ν
 679 is the goal-conditioned target distribution ρ_g which is a Dirac at the given goal state. Similarly, the
 680 distribution to be transported μ is the agent’s goal-conditioned state distribution ρ_π .

681 The Wasserstein-1 distance of an agent executing policy π to the goal s_g can be expressed in a fairly
 682 straightforward manner as:

$$W_1(\rho_\pi, \rho_g) = \sum_{s \in \mathcal{S}} \rho_\pi(s | s_g) d(s, s_g) \quad (14)$$

683 The above is a simplification of Equation 1 where $p = 1$ and the joint distribution is easy to specify
 684 since the target distribution ρ_g is a Dirac distribution.

685 C Lipschitz constant of Potential function

686 For a given goal s_g and all states $s_0 \in \mathcal{S}$, recall that function f is L -Lipschitz if it follows the
 687 Lipschitz condition as follows.

$$f(s_g) - f(s_0) \leq L d_T^\pi(s_0, s_g) \quad \forall s_0 \in \mathcal{S} \quad (15)$$

688 **Proposition 4.** *If transitions from the agent policy π are guaranteed to arrive at the goal in finite*
 689 *time and f is L -bounded in expected transitions, i.e.,*

$$\sup_{s \in \mathcal{S}} \mathbb{E}_{s' \sim \pi, P} [f(s') - f(s)] \leq L,$$

690 *then f is L -Lipschitz.*

691 *Proof.* Since $f(s_g) - f(s_0)$ is a scalar quantity, we may write $f(s_g) - f(s_0) = \mathbb{E}_{\pi, P} [f(s_g) - f(s_0)]$.
 692 Using this fact and that $P(T(s_0) < \infty) = 1$ where $T(s_0) = T^\pi(s_g | \pi, s_0)$ for notation simplicity,

693 the LHS of the expression above becomes a telescopic sum

$$\begin{aligned} f(s_g) - f(s_0) &= \mathbb{E}_{\pi, P} [f(s_g) - f(s_0)] \\ &= \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{T(s_0)-1} f(s_{t+1}) - f(s_t) \right]. \end{aligned}$$

694 Now let us assume that for all transitions (s, a, s') , $\mathbb{E}[f(s') - f(s)] \leq L$. Then

$$\begin{aligned} \mathbb{E}_{\pi, P} \left[\sum_{t=0}^{T(s_0)-1} f(s_{t+1}) - f(s_t) \right] &= \mathbb{E}_{T(s_0)} \left[\mathbb{E}_{\pi, P} \left[\sum_{t=0}^{T(s_0)-1} f(s_{t+1}) - f(s_t) \middle| T(s_0) \right] \right] \\ &\leq \mathbb{E}_{T(s_0)} \left[\sum_{t=0}^{T(s_0)-1} L \right] \\ &= L \mathbb{E}_{T(s_0)} [T(s_0)] \\ &= L d_T^\pi(s_0, s_g), \end{aligned}$$

695 showing that $f(s_g) - f(s_0) \leq L d_T^\pi(s_0, s_g)$. \square

696 D Proofs of Claims

697 The Bellman optimality condition gives us the following optimal distance to goal:

$$d_T^\diamond(s, s_g) = \begin{cases} 0 & \text{if } s = s_g \\ 1 + \min_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} P(s'|s, a, s_g) d_T^\diamond(s', s_g) & \text{otherwise} \end{cases} \quad (16)$$

698 **Proposition 1.** A lower bound on the value of any state under a policy π can be expressed in terms
699 of the time-step distance from that state to the goal: $V(s_0|s_g) \geq \gamma d_T^\pi(s_0, s_g)$.

Proof.

$$V^\pi(s|s_g) = \mathbb{E} \left[\gamma^{T(s_g|\pi, s)} \right] \geq \gamma d_T^\pi(s, s_g) \quad \forall s \in \mathcal{S}$$

700 where the inequality follows as a consequence of Jensen's inequality and the convex nature of the
701 value function. \square

702 **Proposition 2.** If the transition dynamics are deterministic, the policy that maximizes expected return
703 is the policy that minimizes the time-step metric ($\pi^* = \pi^\diamond$).

704 *Proof.* Consider the value of a state s given goal s_g . If the transitions are deterministic and the agent
705 policy π is deterministic (as is the case for the optimal policy), then the time to reach the goal satisfies
706 $\text{Var}(T(s_g|\pi, s)) = 0$, implying that Δ_{Jensen} vanishes and therefore

$$V^\pi(s|s_g) = \gamma^{d_T^\pi(s, s_g)}.$$

707 Since $\gamma \in [0, 1)$, V^π is monotonically decreasing with d_T^π

$$\arg \max_{\pi} V^\pi(s|s_g) = \arg \min_{\pi} d_T^\pi(s, s_g) \quad \forall s \in \mathcal{S}$$

708 That is, in the deterministic transition dynamics scenario, $\pi^* = \pi^\diamond$. \square

709 **Proposition 3.** For a given policy π , the Wasserstein distance of the state visitation measure of that
710 policy from the goal state distribution ρ_g under the ground metric d_T^π can be written as

$$W_1^\pi(\rho_\pi, \rho_g) = \mathbb{E}_{s_0 \sim \rho_0} \left[h(d_T^\pi(s_0, s_g)) + \frac{\gamma}{1 - \gamma} (\Delta_{\text{Jensen}}^\pi(s_0) - 1) \right] \quad (6)$$

711 where h is an increasing function of d_T^π .

712 *Proof.* The first step of the proof is to obtain an analytical expression for the the expected distance to
713 the goal after t steps as a function of the expected distance at $t = 0$. To reduce the notation burden,
714 denote $T(s_0) = T(s_g | \pi, s_0)$ and let $s_t(s_0)$ be the state after t steps conditional on some starting state
715 s_0 where actions are taken according to π . We have excluded s_g and π from the notation since they
716 are fixed for the purpose of this proposition. Using the law of total expectation we have that for every
717 initial s_0

$$\mathbb{E}_{s_t}[d(s_t(s_0), s_g)] = \mathbb{E}_{T(s_0)}[\mathbb{E}_{s_t}[d(s_t(s_0), s_g) | T(s_0)]] = \mathbb{E}_{T(s_0)}[\max(T(s_0) - t, 0)],$$

718 Now, by expanding the definition of $\rho_\pi(s | s_g)$ in equation 5, exchanging the order of summation,
719 and using the previous equation we may write

$$\begin{aligned} W_1^\pi(\rho_\pi, \rho_g) &= \sum_{s \in \mathcal{S}} \sum_{t=0}^{\infty} (1 - \gamma) \gamma^t \mathbb{E}_{s_0}[P(s_t = s | \pi, s_g)] d_T^\pi(s, s_g) \\ &= \mathbb{E}_{s_0} \left[(1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s_t}[d(s_t(s_0), s_g) | s_0] \right] \\ &= \mathbb{E}_{s_0} \left[\mathbb{E}_{T(s_0)} \left[(1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \max(T(s_0) - t, 0) \middle| s_0 \right] \right] \end{aligned}$$

720 Standard but tedious algebraic manipulations given in Lemma 1 in the Appendix show that

$$\sum_{t=0}^{\infty} (1 - \gamma) \gamma^t \max(T(s_0) - t, 0) = T(s_0) - \frac{\gamma}{1 - \gamma} (1 - \gamma^{T(s_0)}).$$

721 Combining the two identities above we arrive at

$$\begin{aligned} W_1^\pi(\rho_\pi, \rho_g) &= \mathbb{E}_{s_0} \left[\mathbb{E}_{T(s_0)} \left[T(s_0) - \frac{\gamma}{1 - \gamma} (1 - \gamma^{T(s_0)}) \middle| s_0 \right] \right] \\ &= \mathbb{E}_{s_0} \left[d(s_0, s_g) - \frac{\gamma}{1 - \gamma} (1 - \mathbb{E}[\gamma^{T(s_0)} | s_0]) \right] \\ &= \mathbb{E}_{s_0} \left[d(s_0, s_g) + \frac{\gamma}{1 - \gamma} \gamma^{d(s_0, s_g)} - \frac{\gamma}{1 - \gamma} (1 - \mathbb{E}[\gamma^{T(s_0)} | s_0] + \gamma^{d(s_0, s_g)}) \right] \\ &= \mathbb{E}_{s_0} \left[d(s_0, s_g) + \frac{\gamma}{1 - \gamma} \gamma^{d(s_0, s_g)} + \frac{\gamma}{1 - \gamma} (\Delta_{\text{Jensen}}^\pi(s_0) - 1) \right]. \end{aligned} \tag{17}$$

722 To finalize the proof, we only need to show that the function $h(\mu) = \mu + \frac{\gamma}{1 - \gamma} \gamma^\mu$ is monotonically
723 increasing for every $\gamma \in [0, 1)$. This is a standard calculus exercise that we show in Lemma 2 in
724 Appendix E. \square

725 **Theorem 1.** *If the transition dynamics are deterministic, the policy that minimizes the Wasserstein*
726 *distance over the time-step metrics in a goal-conditioned MDP (see equation 5) is the optimal policy.*

727 *Proof.* Proposition 2 shows that the Jensen gap vanishes for the optimal policy of an MDP with
728 deterministic transitions and that it minimizes the expected distance from start for all initial states.
729 Proposition 3 on the other hand, implies that when the Jensen gap vanishes, the Wasserstein distance
730 is monotonically increasing in the expected distance from the start. Together, the two propositions
731 show that π^* minimizes the Wasserstein distance. \square

732 E Auxiliary results for Proposition 3

733 **Lemma 1.** *Let T be a positive integer. Then*

$$\sum_{t=0}^{\infty} (1 - \gamma) \gamma^t \max(T - t, 0) = T - \frac{\gamma}{1 - \gamma} (1 - \gamma^T).$$

Algorithm 1: AIM + HER

Input: Agent policy π_θ , discriminator f_ϕ , environment env , number of Epochs N , number of time-steps per epoch K , policy update period k , discriminator update period m , episode length T , replay buffer (for HER), smaller replay buffer (for discriminator)

```
1 Initialize discriminator parameters  $\phi$ ;  
2 Initialize policy parameters  $\theta$ ;  
3 for  $n = 0, 1, \dots, N - 1$  do  
4    $t = 0$ ;  
5    $goal\_reached = \text{True}$ ;  
6   while  $t < K$  do  
7     if  $goal\_reached$  or  $episode\_over$  then  
8       Sample goal  $s_g \sim \sigma(\mathcal{G})$ ;  
9       Sample start state  $s \sim \rho_0(\mathcal{S})$ ;  
10       $goal\_reached = \text{False}$ ;  
11       $episode\_over = \text{False}$ ;  
12       $t_{start} = K$ ;  
13     end  
14     Sample action  $a \sim \pi_\theta(\cdot | s, s_g)$ ;  
15      $s' = env.step(a)$ ;  
16     if  $s' = s_g$  then  
17       |  $goal\_reached = \text{True}$ ;  
18     end  
19     // end episode if goal not reached in  $T$  steps  
20     if  $t - t_{start} = T$  then  
21       |  $episode\_over = \text{True}$ ;  
22     end  
23     Add  $(s, a, s', s_g, goal\_reached)$  to replay buffer and smaller replay buffer;  
24     if  $goal\_reached$  or  $episode\_over$  then  
25       | Add hindsight goals to both buffers;  
26     end  
27     // Update policy parameters  $\theta$  every  $k$  steps  
28     if  $t \% k = 0$  then  
29       | Sample tuples  $(s, a, s', s_g, goal\_reached)$  from replay buffer;  
30       | Get intrinsic reward (Equation 9);  
31       | Update policy parameters  $\theta$  using any off-policy learning algorithm;  
32     end  
33     // Update discriminator parameters  $\phi$  every  $m$  steps  
34     if  $t \% m = 0$  then  
35       | Sample tuples  $(s, a, s', s_g, goal\_reached)$  from smaller replay buffer;  
36       | Update discriminator parameters  $\phi$  using Equation 8;  
37     end  
38      $t = t + 1$ ;  
39   end  
40   Evaluate agent policy;  
41 end
```

734 *Proof.* Direct computation gives

$$\begin{aligned} (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \max(T - t, 0) &= (1 - \gamma) \sum_{t=0}^{T-1} \gamma^t (T - t) \\ &= (1 - \gamma) T \sum_{t=0}^{T-1} \gamma^t - (1 - \gamma) \sum_{t=0}^{T-1} t \gamma^t \end{aligned}$$

735 We will now simplify the two terms of the last expression. For the first one, have

$$(1 - \gamma)T \sum_{t=0}^{T-1} \gamma^t = (1 - \gamma)T \frac{1 - \gamma^T}{1 - \gamma} = T - T\gamma^T.$$

736 For the second one, the computations are a bit more involved

$$\begin{aligned} (1 - \gamma) \sum_{t=0}^{T-1} t\gamma^t &= (1 - \gamma)\gamma \sum_{t=1}^{T-1} t\gamma^{t-1} \\ &= (1 - \gamma) \sum_{t=1}^{T-1} \gamma \frac{d}{d\gamma} \gamma^t \\ &= \gamma(1 - \gamma) \frac{d}{d\gamma} \sum_{t=0}^{T-1} \gamma^t \\ &= \gamma(1 - \gamma) \frac{d}{d\gamma} \frac{1 - \gamma^T}{1 - \gamma} \\ &= \frac{\gamma}{(1 - \gamma)} (-T\gamma^{T-1}(1 - \gamma) + (1 - \gamma^T)) = -T\gamma^T + \frac{\gamma}{(1 - \gamma)}(1 - \gamma^T). \end{aligned}$$

737 When combining the two simplified expressions the terms with $T\gamma^T$ will cancel out, yielding the
738 desired expression. \square

739 **Lemma 2.** *The function $h_\gamma(\mu) = \mu + \frac{\gamma}{1-\gamma}\gamma^\mu$ is monotonically increasing for every $\gamma \in [0, 1)$.*

740 *Proof.* We must show that $\frac{d}{d\mu}h_\gamma(\mu) > 0$ for every $\gamma \in [0, 1)$ and every $\mu > 0$. Computing the
741 derivative directly we obtain

$$\frac{d}{d\mu}h_\gamma(\mu) = 1 + \frac{\log(\gamma)\gamma^{\mu+1}}{1 - \gamma}.$$

742 Thus, it will suffice to show that the second term above is greater than -1. For this purpose, first note
743 that $\log(\gamma)\gamma^{\mu+1} > \log(\gamma)$ since $\gamma < 1$. Now, we use the fact that $\log(\gamma) < 1 - \gamma$ for $\gamma < 1$. This
744 can be verified noting that $1 - \gamma$ is the tangent line to the concave curve $\log(\gamma)$ and the curves meet
745 at $\gamma = 1$. And therefore $\log(\gamma)/(1 - \gamma) > -1$. Putting these observation together,

$$\frac{d}{d\mu}h_\gamma(\mu) = 1 + \frac{\log(\gamma)\gamma^{\mu+1}}{1 - \gamma} > 1 + \frac{\log(\gamma)}{1 - \gamma} > 1 - 1 = 0,$$

746 concluding the proof. \square

747 F Grid World Experiments

748 **Basic experiment** The environment is a 10×10 grid with 4 discrete actions that take the agent in
749 the 4 cardinal directions unless blocked by a wall or the edge of the grid. The agent policy is learned
750 using maximum entropy Q-learning [30], with an entropy coefficient of 0.1 and a discount factor of
751 $\gamma = 0.99$. We do not use hindsight goals for this experiment, and use a single buffer with size 5000
752 for both the policy as well as the discriminator training. The results are discussed in the main text.
753 The compute used to conduct these experiments was a personal laptop with an Intel i7 Processor and
754 16 GB of RAM.

755 **Additional experiments** We conducted variations from the basic experiment in the grid world to
756 show that AIM and its novel regularization can learn a reward function which guides the agent to the
757 goal even in the presence of stochastic transitions as well as transitions where the state features vary
758 wildly from one step to the next.

759 First, we evaluate AIM’s ability to learn in the presence of stochastic and asymmetric transitions in a
760 windy version (Figure 4a) of the above grid world. Transitions in the last six columns of the grid are
761 affected by a wind blowing from the top. Actions that try to move upwards only succeed 60% of the

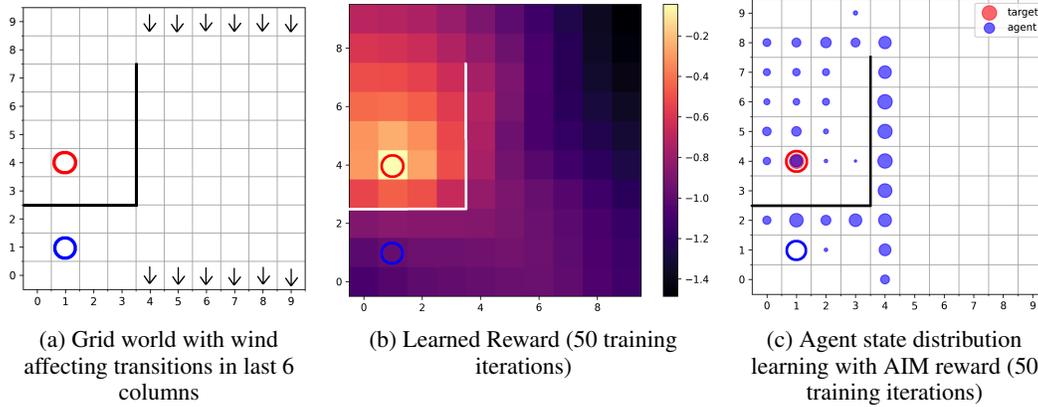


Figure 4: Windy grid world (Figure 4a) experiments. The columns with arrows at the top and bottom have stochastic and asymmetric transitions induced by wind blowing from the top. Learned reward function (Figure 4b). Reward at each state of the grid world after training for 50 iterations with AIM. Hollow red circle indicates the goal state. White lines indicate the walls the agent cannot transition through. The agent’s state visitation (Figure 4c): The hollow blue circle indicates agent’s start state. The hollow red circle is the goal. Blue bubbles indicate relative time the agent’s policy causes it to spend in respective states. Black lines indicate walls.

762 time, and actions attempting to move sideways cause a transition diagonally downwards 40% of the
 763 time. Movements downwards are unaffected. The rest of the experiment is carried out in the same
 764 way as above, but with 128 hidden units in the hidden layer of the agent’s Q function approximator
 765 (the reward function architecture is unchanged from the previous experiment). In Figure 4 we see
 766 that AIM learns a reward function that is still useful and interpretable, and leads to a policy that can
 767 confidently reach the goal, regardless of these stochastic and asymmetric transitions. Notice the effect
 768 of the stochastic transitions in the increased visitation in the sub-optimal states in the bottom two
 769 rows of column number 4.

770 The next experiment tests what happens when the transition function causes the agent to jump between
 771 states where the state features vary sharply. As an example consider a toroidal grid world, where
 772 if an agent steps off one side of the grid it is transported to the other side. The distance function
 773 here should be smooth across such transitions, but might be hampered by the sharp change in input
 774 features. In Figure 5 we see show the policy and reward for a 10×10 toroidal grid world with start
 775 state at (2, 2) and goal at (7, 7). Transitions are deterministic but wrap around the edges of the grid
 776 as described above: a **down** action in row 0 will transport the agent to the same column but row 9.
 777 The start and the goal state are set up so that there are multiple optimal paths to the goal. The entropy
 778 maximizing soft Q-learning algorithm should take these paths with almost equal probability. From
 779 Figure 5 it is evident that AIM learns a reward function that is smooth across the actual transitions in
 780 the environment and allows the agent to learn a Q-function that places near equal mass on multiple
 781 trajectories.

782 G Statistical Analysis of the Results on Fetch Robot Tasks

783 To compare the performance of each method with statistical rigor, we used a repeated measures
 784 ANOVA design for binary observation where an observation is successful if an agent reaches the goal
 785 within an episode. We then conducted a Tukey test to compare the effects of each method, i.e., the
 786 estimated odds of reaching the goal given the algorithm. The goal of the statistical analysis presented
 787 here is twofold

- 788 1. Separate the uncertainty on the performance of each method from the variation due to
 789 random seeds.
- 790 2. Adjust the probability of making a false discovery due to multiple comparisons. This extra
 791 step is necessary to avoid detecting a large fraction of falsely “significant” differences since
 792 typical tests are designed to control the error rate of only one experiment.

793 The data for statistical analysis comes from $N_{\text{episodes}} = 100$ evaluation episodes per each one of
 794 $N_{\text{seeds}} = 6$ seeds. For all environments but FetchReach, these data is collected after 1 million
 795 environment interactions; and for FetchReach it is taken after 2000 interactions.

796 The repeated measures ANOVA design is formulated as a mixed effects generalized linear model and
 797 fitted separately for each one of the four environments

$$\begin{aligned}
 y_{ijk} &\stackrel{\text{iid}}{\sim} \text{Bernoulli}(p_{ij}), & k &\in \{1, \dots, N_{\text{episodes}}\} \\
 \text{logit}(p_{ij}) &= r_{\text{seed}_i} + \beta_{\text{algorithm}_j}, & i &\in \{1, \dots, N_{\text{seeds}}\}, j \in \{1, \dots, N_{\text{algorithms}}\} \\
 r_{\text{seed}_i} &\stackrel{\text{iid}}{\sim} \text{Normal}(0, \sigma^2)
 \end{aligned}$$

798 The variation due to the seed effects is measured by σ^2 , whereas the uncertainty about the odds of
 799 reaching the goal using each algorithm is measured by the standard errors of the coefficients $\beta_{\text{algorithm}_j}$.
 800 The Tukey test evaluates all null hypotheses $H_0: \beta_{\text{algorithm}_j} = \beta_{\text{algorithm}_{j'}}$ for all combinations of j, j' .
 801 To adjust for multiple comparisons each Tukey tests uses the Holm method. Since we are also doing
 802 a Tukey test for each environment, we further apply a Bonferroni adjustment with a factor of four.
 803 These types of adjustments are fairly common for dealing with multiple comparison in the literature
 804 of experimental design; the interested reader may consult [45].

805 The results, shown in Table I, signal strong statistical evidence of the improvements from using
 806 the AIM learned rewards. In three of the four environments AIM and AIM+ R have similar odds of
 807 reaching the goal as the dense shaped reward (H_0 is not rejected,) and in all four environments AIM
 808 and AIM+ R have statistically significant higher odds of reaching the goal than the sparse reward (H_0
 809 is rejected and β is higher.)

Contrast	Slide	Push	PickAndPlace	Reach
$\beta_{\text{AIM+R}} - \beta_{\text{HER+dense}}$	0.34 (0.14)	-1.74 (0.77)	-0.10 (0.45)	*-3.43 (0.34)
$\beta_{\text{AIM}} - \beta_{\text{HER+dense}}$	0.21 (0.14)	-2.19 (0.75)	*-1.50 (0.37)	*-5.01 (0.35)
$\beta_{\text{AIM+R}} - \beta_{\text{HER+sparse}}$	*0.69 (0.13)	*5.32 (0.35)	*4.71 (0.33)	*4.75 (0.25)
$\beta_{\text{AIM}} - \beta_{\text{HER+sparse}}$	*0.57 (0.13)	*4.86 (0.30)	*3.31 (0.19)	*3.17 (0.24)

Table 1: Results of the Tukey test on the evaluation of Fetch tasks. The table entries are log odds ratios with standard deviations shown in parentheses. Positive values mean that AIM or AIM+R perform better than the method with negative sign in the contrast and viceversa. Asterisks mark statistical significance at 95%. If there is no asterisk, then H_0 is not rejected in which case the differences could be due to random chance.

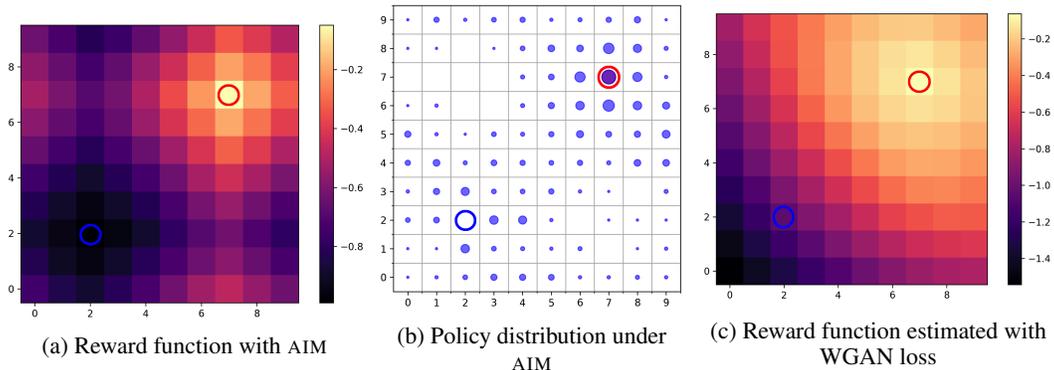


Figure 5: The reward function (Figure 5a) learned with AIM and subsequent policy distribution (Figure 5b) in a toroidal grid world, where an agent can transition from one edge of the grid across to the other. The hollow blue circle denotes the start state and the hollow red circle is the goal state. The reward function respects the sharp transitions from one end of the grid to the other. Conversely, if the reward function is learned using the WGAN objective [29] (Figure 5c), it does not respect the environment dynamics.

810 **H Details of Experiments on Fetch Robot**

811 The Fetch robot domain in OpenAI gym has four tasks available for testing. They are named Reach,
 812 Push, Slide, and Pick And Place. The Reach task is the simplest, with the goal being the 3-d
 813 coordinates where the end effector of the robot arm must be moved to. The Push task requires pushing
 814 an object from its current position on the table to the given target position somewhere else on the
 815 table. Slide is similar to Push, except the coefficient of friction on the table is reduced (causing
 816 pushed objects to slide) and the potential targets are over a larger area, meaning that the robot needs
 817 to learn to hit objects towards the goal with the right amount of force. Finally, Pick And Place is the
 818 task where the robot actuates it’s gripper, picks up an object from its current position on the table and
 819 moves it through space to a given target position that could be at some height above the table. The
 820 goal space for the final three tasks are the required position of the object, and the goal the current
 821 state represents is the current position of that object.

822 Next, we note the hyperparameters used for various baselines as well as our implementation. The
 823 names of the hyperparameters are as specified in the stable baselines repository and used in the RL
 824 Zoo [54] codebase which we use for running experiments. Both the stable baselines repository and
 825 RL Zoo are available under the MIT license. These experiments were run on a compute cluster with
 826 each experiment assigned an Nvidia Titan V GPU, a single CPU and 12 GB of RAM. Each run of the
 827 TD3 baseline HER + R or HER + dense required 18 hours to execute, and each run which included
 828 AIMrequired 24 hours to complete execution.

829 TD3 [23], like its predecessor DDPG [43], suffers from the policy saturating to extremes of its
 830 parameterization. Hausknecht and Stone [32] have suggested various techniques to mitigate such
 831 saturation. We use a quadratic penalization for actions that exceed 80% of the extreme value at
 832 either end, which is sufficient to not hurt learning and prevent saturation. Assuming the policy
 833 network predicts values between -1 and 1 (as is the case when using the tanh activation function),
 834 the regularization loss is:

$$L_a = \frac{1}{N} \sum_{i=1}^N [\max(|\pi_\theta(s_i)| - 0.8, 0)]^2$$

835 where N is the mini-batch size and s_i is the state for the i^{th} transition in the batch.

836 The other modification made to the stable baselines code is to use the Huber loss instead of the
 837 squared loss for Q-learning.

838 For evaluation, in the Reach domain the agent policy is evaluated for 100 episodes every 2000 steps.
 839 For the other three domains, the experiment is run for 1 million timesteps, and evaluated at every
 840 20,000 steps for 100 episodes.

841 **H.1 TD3 and HER (R + HER)**

Hyperparameter	Value
n_sampled_goal	4
goal_selection_strategy	future
buffer_size	10^6
batch_size	256
γ (discount factor)	0.95
random_exploration	0.3
target_policy_noise	0.2
learning_rate	1^{-3}
noise_type	normal
noise_std	0.2
MLP size of agent policy and Q function	[256, 256, 256]
learning_starts	1000
train_freq	10
gradient_steps	10
τ (target policy update rate)	0.05

843 **H.2 Dense reward TD3 and HER (dense + HER)**

Hyperparameter	Value
n_sampled_goal	4
goal_selection_strategy	future
buffer_size	10^6
batch_size	256
γ (discount factor)	0.95
random_exploration	0.3
target_policy_noise	0.2
844 learning_rate	1^{-3}
noise_type	normal
noise_std	0.2
MLP size of agent policy and Q function	[256, 256, 256]
learning_starts	1000
train_freq	100
gradient_steps	200
policy_delay	5
τ (target policy update rate)	0.05

845 **H.3 TD3 and HER with AIM (AIM + HER) and (AIM + R + HER)**

Hyperparameter	Value
n_sampled_goal	4
goal_selection_strategy	future
buffer_size	10^6
batch_size	256
γ (discount factor)	0.9
random_exploration	0.3
target_policy_noise	0.2
846 learning_rate	1^{-3}
noise_type	normal
noise_std	0.2
MLP size of agent policy and Q function	[256, 256, 256]
learning_starts	1000
train_freq	100
gradient_steps	200
disc_train_freq	100
disc_steps	20
τ (target policy update rate)	0.1