

1 We thank the reviewers for their valuable comments and suggestions, which will help us to improve the paper. We first
 2 respond to **R1: 1)** While “twin-delay” is well-established to help estimate $Q(s, \mathbf{a})$, IDAC advances it to a distributional
 3 RL setting for continuous actions to estimate the distribution of the discounted cumulative return $Z(s, \mathbf{a})$, whose
 4 expectation is $Q(s, \mathbf{a})$. While the usual “twin-delay” only involves minimization of two scalars, the one in IDAC
 5 is distinct in involving *element-wise* minimization of two *sorted* vectors, whose elements are *iid* sampled from two
 6 different DGNs. Besides, we introduce SIA and an asymptotic lower bound for entropy estimation. Both new techniques
 7 are helpful for IDAC to outperform the SOTA algorithms (please refer to Fig. 3). **2)** We will add the suggested ablation
 8 studies. **3)** Similar to Fig. 2(a)(b) illustrating SIA at an early training stage, we have examined SIA at a late stage and
 9 found that its differences from a diagonal Gaussian remain evident, as verified with both a normality test on the marginal
 10 of each dimension of the SIA policy and the Pearson correlation test on many randomly selected action dimension pairs.
 11 In addition, please see a related response to R3 (Lines 28-39) that strengthens the claim that SIA can represent more
 12 complicated policy distributions. These details will be added. **4)** As suggested, we will draw more explicit distinctions,
 13 such as the unique operation of sorting followed by element-wise minimization, and expand related work.

14 **R2: 1)** We will clarify sorting is applied to each individual vector. We denote \vec{x}_i , which is a scalar, as the i -th element
 15 after sorting vector \mathbf{x} . **2)** We will better discuss related work to help avoid confusions and highlight our contributions.

16 **R3: 1)** We are puzzled why IDAC is considered by R3 to be not sufficiently compared against other similar approaches;
 17 we would appreciate R3 pointing out these missing baselines. First, from the actor-critic perspective, we have already
 18 compared IDAC to SOTAs: PPO, SAC, & TD3. Second, from the distributional RL perspective, while there exist
 19 well-known algorithms (e.g., C51, QR-DQN, and IQN) for discrete actions, we find D4PG to be the only published
 20 one designed for continuous controls. Note we did not add a direct comparison to D4PG (as well as SDPG) as its
 21 implementation has notable differences from the Stable-baselines used in this paper. In particular, D4PG uses distributed
 22 agents to do distributed sampling for the replay buffer, allowing an implied advantage in observing more state-action
 23 pairs given the same number of policy gradient update steps. This is the main reason that we have adapted SDGP,
 24 which improves over D4PG, into Stable-baselines and used it for ablation study. To help address R3’s concern, we
 25 will add these SDGP results (clearly worse than IDAC) into Fig. 1. Moreover, we have run the original D4PG code
 26 in hope to further eliminate this concern, and found that it consistently underperforms IDAC given the same number
 27 of policy gradient update steps. E.g., the max average returns of [D4PG, IDAC] after 10^6 policy gradient steps are
 28 $[9776 \pm 739, 12222 \pm 157]$ on HalfCheetah-v2, and $[4742 \pm 1320, 5386 \pm 335]$ on Walker2d-v2. **2) Verifying that**
 29 **SIA would improve exploration:** We note good exploration by SIA can be implied from Fig. 3, by better empirical
 30 performance; as the reward mechanism in continuous control RL tasks is complicated, such empirical comparison is
 31 a common way to indirectly evaluate whether better explorations have been achieved for these tasks (e.g. Sec. 4.4
 32 of Hong et al. [2018], arXiv:1802.04564). To more directly verify that SIA does improve exploration, we mimic
 33 Sec. 5.1 of Haarnoja et al. [2017], arXiv:1702.08165 to introduce a Multigoal Environment with four equally-spaced
 34 and well-separated destinations in a 2D map that provide large rewards of [200, 200, 200, 600], with spatial location-
 35 dependent negative rewards elsewhere. This task requires extensive exploration to reach the optimal solution. We
 36 combine a Gaussian actor or a SIA with an Actor-Critic (A2C) algorithm; each algorithm is trained with 10 independent
 37 runs (10^5 episodes, one evaluation). In a single run, Gaussian often only explores no more than two (sub-optimal)
 38 destinations, achieving average cumulative reward as 123.17 ± 156.86 , while SIA in general successfully explores all
 39 four destinations, achieving 336.84 ± 26.9 . These details verify that SIA does help exploration. **3)** We’d like to point
 40 out that IDAC uses Huber loss, and it is unclear whether additional theoretical justifications are needed to combine
 41 Huber loss with empirical Wasserstein distance. As quantile regression Huber loss is successfully used by QR-DQN for
 42 discrete actions, avoiding the potential issue of having biased gradients, we made the safe choice of using the same loss
 43 in IDAC. We agree it would be an interesting future work to investigate the use of empirical Wasserstein distance (with
 44 Huber loss) under IDAC. **4)** We will add more textual descriptions to improve our pseudo-code in the Appendix.

45 **R4: 1)** We integrate DGN and SIA to allow them to help each other: (i) Modeling G exploits distributional information
 46 to help better estimate its mean Q (note C51, which outperforms DQN by exploiting distributional information, also
 47 conducts its argmax operation on Q); (ii) A more flexible policy may become more necessary given a better estimated
 48 Q . We have followed your suggestions to conducted additional ablation studies, which show removing either SIA or
 49 DGN from IDAC in general negatively impacts its performance. These results will be added into revision. **2)** Quantile
 50 regression loss can be applied to both explicit (as in QR-DQN) and implicit distribution (as in DGN). We favor DGN
 51 because it not only is more flexible, but also avoids a potential pitfall: QR-DQN fixes the quantile locations and feeds
 52 each (s, \mathbf{a}) to a deep neural network (NN) to estimate their corresponding values; A clear concern, however, is that
 53 given input (s, \mathbf{a}) , the NN output at a designated higher quantile is not guaranteed to be larger than that at a lower
 54 quantile. By contrast, there is no such concern in DGN, which simply sorts its *iid* sampled values. **3)** We clarify that
 55 the SAC results were obtained by using the modern version [41] with the learned temperature. **4)** In Table 1 caption, we
 56 will make the change to use “max average returns”, obtained by taking the maximal of averaged evaluation score (over
 57 4 random seeds). **5)** The normalizing integral does not contribute to the gradient of θ and hence omitted. In Eq(15), it is
 58 $\pi_\theta(\mathbf{a}^{(1)}|s, \xi^{(\ell)})$ rather than $\pi_\theta(\mathbf{a}^{(\ell)}|s, \xi^{(\ell)})$, as indicated in Eq(12) multiple auxiliary variables ξ are needed for each \mathbf{a} .