Figure A    Figure B    Figure C    Figure D

1  We thank the reviewers for insightful comments. We are improving the paper by incorporating the reviewers' suggestions.

2  **Response to reviewer #1:**

3  **Why we present two SBR instances (Table 2)**  The first $d_\phi^2$ requires sampling from a uniform distribution to approx-

4  imate $f$, which leads to a large variance of gradient and then hurts the performance (see Figure A). Therefore, we

5  propose $\tilde{d}_\phi^2$ to reduce the variance, which provides a closed-form expression of $f$ and does not require sampling. We

6  will add more discussions about these two instances in the final version, if accepted.

7  **Discussion about Theorem 5**  Theorem 5 proposes a bound on the performance gap between the regularized optimal

8  policy $\pi_\alpha^*$ and the original optimal policy $\pi^*$. This theorem shows that the gap depends on $\Delta = \phi(U^2) - \phi(0)$. Though

9  larger $\Delta$ can encourage actions away from each other better, it may lead to worse performance of $\pi_\alpha^*$. We will add the

10  discussion about Theorem 5 in the final version, if accepted.

11  **Performance difference between SAC and our method (ACED)**  When we use Gaussian policies, the hyperparam-

12  eters of ACED (Table 6 in Appendix) are the same as that of SAC except for the target value $T$. Therefore, the

13  performance difference mainly comes from different regularization.

14  **More examples that demonstrate the benefits of ACED**  To show the improved efficiency of ACED, we evaluate

15  ACED ($N = 2$) against SAC with an ensemble of policies. The ensemble size is 5, and each policy outputs a Gaussian

16  distribution. Experiments show that 1000 updates in ACED cost 21.4s while those in SAC cost 33.7s. The core reason

17  is that ACED does not need to compute probability density, which requires the forward prediction of all networks. We

18  will provide results in details in the final version, if accepted.

19  **Response to reviewer #2:**

20  **Compare the two SBR instances**  In practice, $\tilde{d}_\phi^2$ is better than $d_\phi^2$. Please see the response to reviewer #1 for details.

21  **Histograms of entropies of ACED and SAC**  In Figure B, we provide the histograms of entropies of SAC and ACED

22  during training. It shows that the difference of stochasticity exists through almost the whole training procedure.

23  **Response to reviewer #3:**

24  **The meaning of "we did not tune the hyperparameters"**  For SAC and TD3, we use the hyperparameters provided

25  by their authors. Therefore, we did not tune them again. For our algorithm, we use the same hyperparameters (see

26  Table 6 in Appendix) as SAC, if possible. We did not tune $N$ and $T$ for the results in Section 5.1.

27  **Examples of complex policies**  We provided evaluation with normalizing flow policies as an example in Section 5.2.

28  Moreover, expressing the policies by a noisy network [10] requires an additional classification network to estimate

29  entropy [44]. We will provide more examples in the final version, if accepted.

30  **Response to reviewer #4:**

31  **Discussion about Theorem 5**  Please refer to the response to reviewer #1.

32  **Should the community switch from SAC to ACED?**  If computing the probability density of policies is time-

33  consuming or infeasible, the answer is "yes". For example, when parameterizing policies by noisy networks, ACED is a

34  better choice than SAC. Otherwise, the answer depends on the performance of SAC and ACED.

35  **How sensitive ACED is to hyperparameters?**  ACED is insensitive to the sample number $N$ and the target value $T$.

36  The sensitivity analysis for $N$ can be found in Section 5.3 and Appendix C.5. Figure C shows the sensitivity analysis

37  for $T$. Here, we set $T$ as $\mathcal{F}(\mathcal{N}(\mathbf{0}, \lambda\mathbf{I}))$, where $\mathcal{N}(\mathbf{0}, \lambda\mathbf{I})$ is a Gaussian distribution with the covariance matrix $\lambda\mathbf{I}$.

38  **Connections with related work and our novelty contributions**  Most existing regularization [18,23,51] takes the

39  form of $\mathbb{E}_{a\sim\pi(\cdot|s)}[f(\pi(a|s))]$. We propose a novel regularization form (Equation 4) to encourage stochasticity. Due to

40  the different forms, previous regularization often requires computing probability density to estimate entropy but our

41  regularization does not. Moreover, unlike previous regularization, our regularization considers the distances between

42  actions and thus incorporates geometric information.

43  **What is Gini mean difference?**  Gini mean difference is a measure of statistical dispersion. Given a distribution $D$, it

44  is defined by $\mathbb{E}|X_1 - X_2|$, where $X_1$ and $X_2$ are sampled from $D$ independently.

45  **What is the difference between Figures 1(a) and 1(d)?**  The reward function in Figure 1(a) is unimodal, while the

46  other is trimodal. Figure 1(d) shows that our regularization can lead to a multi-modal distribution.

47  **What does the notation of $\phi^{(n+1)}$ mean on line 66?**  $\phi^{(n+1)}$ denote the $(n + 1)$th order derivative of $\phi$.

48  **Comparison with methods that do not require a specific form of policy:**  We compare ACED with GAC [47] in

49  HalfCheetah-v2 . Figure D shows that our method outperforms GAC. We note that GAC is expensive in computation.

50  GAC takes more than 55h to train the policy with about 0.8 million steps, while ACED takes less about 6h. As the

51  experiments are still running, we will provide more results in the final version, if accepted.