1 We greatly appreciate the reviewers' effort and helpful comments. We are improving the paper by incorporating the
2 reviewers' suggestions.

3 **R#1** 1. A policy has multi-modality means it assigns positive probability mass to optimal and near optimal actions.
4 This concept is more useful when the multiple optimal (or near optimal) actions exist. In this case, a policy with
5 multi-modality is a distribution over these (near) optimal actions and provides information of multiple actions (including
6 the probability weight) in a state, unlike the deterministic one which provides only one action information.

7 2. Based on our results, among the four basic regularizers, cos has the strongest sparsity power (Lines 257-259). This
8 regularizer performs better in scenarios with high-dimensional action space. For example, the performance of cos
9 increases fastest in the environment of Seaquest where $|\mathcal{A}| = 18$ (see (d) in Figure 2). Thanks!

10 **R#2** 1. Line 734: Yes, it's Tanh transformation. The random variable is transformed as $a_t = \tanh(Z_t)$, and the
11 corresponding density relation is $\pi(a_t|s_t) = \mathcal{N}(Z_t)|\det(\frac{da_t}{dZ_t})|^{-1}$, where $\log|\det(\frac{da_t}{dZ_t})| = \sum_{i=1}^{\mathcal{A}} \log(1-\tanh^2(Z_{t,i}))$.
12 We have taken it into account but miss it in the appendix. We will add up the details in its revision.

13 2. Tanh sensible choice?: (i) Indeed, the similar theoretical results can be derived in continuous action space by
14 variational calculus. As $\lambda$ goes zero, the density function reduces to Dirac function centered at $\arg\max_a Q^*(s,a)$.
15 Besides, the density reduces to zero exactly for $a \notin \arg\max_a Q^*(s,a)$ when sparse regularizers are applied. But, in
16 empirical continuous action space setting, it is difficult to find a tractable way to model a function class containing
17 Dirac function and density functions that can be sparse. (You mention that some prior densities can be used to induce
18 sparsity, which is a good advice. We will consider it in the future.) So we just follow the Mujoco experiment setting in
19 (Haarnoja et al. "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor." )
20 to check whether the algorithm still performs well with our proposed regularizers.

21 (ii) In discrete action space, we require the trigonometric family functions are non-negative as the usual entropy (e.g.
22 Shannon entropy), which is lower bounded by 0. The added non-negative term can be viewed as a bonus.

23 3. Line 287: These regularizers have multi-modality (which we have explained in R#1.1) as soft-max policy does. But
24 soft-max policy doesn't have sparsity since it assigns positive probabilities to all actions.

25 4. Lines 123-124: The only consideration is whether the regularizer is differentiable. For ease of analysis, we only
26 consider the case where $\phi(x)$ is fully differentiable. Actually, we can relax (3) in Assumption 1 such that $\phi(x)$ is
27 differential except for finite points. And our theoretical results still hold in this case.

28 5. CVXOPT or GD: It depends on the size of problem. If the problem is simple (small state and action space), CVXOPT
29 converges faster. But if the problem is complicated (e.g., Atari) with large samples, CVXOPT is slow as we have to
30 solve a convex problem for each sample $(s,a,s',r)$ in a batch, which is time-consuming when $|\mathcal{A}|, |\mathcal{S}|$ are large. GD is
31 more scalable by contrast. Thanks!

32 **R#3** Thanks for your careful review. We will incorporate your comments and suggestions into the revision.

33 1. Line 17: Sorry, our phrasing is not careful. Here our point is that greedily solving the Bellman equation is not an easy
34 task in a high-dimension action space or when function approximation (such as neural networks) is used.

35 2. Lines 18-19: Yes, it is possible to break ties randomly. But our point is to emphasize the need of information of
36 multiple actions. Multiple (near) optimal actions provide alternatives when the suggested action is suddenly forbidden.

37 3. Line 21: If the task is path planning where the state is the pair of departure and destination, when the suggested routine
38 is unfortunately congested, an alternative routine could be provided by a multi-modal policy. In this way, we don't need
39 to evoke the computation of new routines. The example aims to shed light on the importance of multi-modality.

40 4. Line 24: Yes, we want the learning problem is computationally inexpensive and the optimal policy has multi-modality.

41 5. Line 27: Yes, it's unclear. We will delete it.

42 6. Line 31-32: Yes, we don't want terrible actions have a chance to be executed. We will modify this sentence.

43 7. Line 49-50: The point is that there exist other regularizers that induces sparsity.

44 8. Line 70: It's a typo. We require $\gamma$ to be non-negative in our proof and our theories still hold when $\gamma = 0$.

45 9. Lines 102-103: Yes, I agree with you. We will add this distinction in the revision.

46 10. Theorem 1: The primal problem (Eqn. 2) can be transformed into a convex problem about the vector $\{\pi(a|s)\}_{s,a}$,
47 so Theorem 1 is necessary and sufficient as it is derived from the KKT conditions.