

1 Thank you for the constructive comments and suggestions. We will incorporate all presentation improvements suggested.

2 **Theoretical results (Reviewer 1):** [ω^\top and \arg_Q cancel each other] By definition (line 143),
 3 $\arg_Q \sup_{a \in \mathcal{A}, \omega' \in \Omega} \omega^\top Q(s, a, \omega') := Q(s, a', \omega'')$, i.e., the \arg_Q operator extracts the Q value that results in the largest
 4 utility using the preferences ω . Therefore, linearizing this Q with the same ω results in exactly the same supremum, i.e.,
 5 $\omega^\top \arg_Q \sup_{a \in \mathcal{A}, \omega' \in \Omega} \omega^\top Q(s, a, \omega') = \omega^\top Q(s, a', \omega'') = \sup_{a \in \mathcal{A}, \omega' \in \Omega} \omega^\top Q(s, a', \omega')$. **Note the supremum is over ω' , not ω .**

6 **[Theorem 1]** Thanks for catching the typo - Q should be Q^* . We realize that the proofs are a bit compressed - we will
 7 update the paper with more detailed derivations for all proofs. Here is Thm. 1 in detail (starting step 2 under line 601):

$$\begin{aligned}
 \omega^\top \mathcal{T}Q^*(s, a, \omega) &= \omega^\top r(s, a) + \gamma \cdot \omega^\top \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a)} \arg_Q \sup_{a' \in \mathcal{A}, \omega' \in \Omega} \omega^\top Q^*(s', a', \omega') \\
 \text{(linearity of exp. \& cancel } \omega^\top \text{ and } \arg_Q) &= \omega^\top r(s, a) + \gamma \cdot \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a)} \sup_{a' \in \mathcal{A}, \omega' \in \Omega} \omega^\top Q^*(s', a', \omega') \\
 \text{(insert eq. (20), def. of } Q^*) &= \omega^\top r(s, a) + \gamma \cdot \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a)} \sup_{a' \in \mathcal{A}, \omega' \in \Omega} \omega^\top \left\{ \arg_Q \sup_{\pi \in \Pi} \omega'^\top \mathbb{E}_{\tau \sim (\mathcal{P}, \pi)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] \right\} \\
 \text{(use def. of } \arg_Q, \text{ explained below)} &= \omega^\top r(s, a) + \gamma \cdot \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a)} \sup_{a' \in \mathcal{A}} \omega^\top \left\{ \arg_Q \sup_{\pi \in \Pi} \omega'^\top \mathbb{E}_{\tau \sim (\mathcal{P}, \pi)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] \right\} \\
 \text{(rearrange expectation and sup)} &= \omega^\top r(s, a) + \gamma \cdot \omega^\top \arg_Q \sup_{\pi \in \Pi} \omega'^\top \mathbb{E}_{s_0 \sim \mathcal{P}(\cdot | s, a)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] \\
 \text{(merge 1st term to sum \& use def. of } Q^* \text{ again)} &= \omega^\top \left\{ \arg_Q \sup_{\pi \in \Pi} \omega'^\top \mathbb{E}_{s_0 \sim \mathcal{P}(\cdot | s, a)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] \right\} = \omega^\top Q^*(s, a, \omega)
 \end{aligned}$$

8 The fourth equation is due to a sandwich inequality, $\omega^\top \arg_Q \sup_{\pi \in \Pi} \omega^\top Q^\pi \leq \sup_{\omega' \in \Omega} \omega^\top \arg_Q \sup_{\pi \in \Pi} \omega'^\top Q^\pi = \omega^\top \arg_Q \sup_{\pi \in \Pi} \omega_*^\top Q^\pi =$
 9 $\omega^\top Q^{\pi^* \omega'_*} \leq \omega^\top \arg_Q \sup_{\pi \in \Pi} \omega^\top Q^\pi$, where ω'_* and $\pi^*_{\omega'_*}$ are preference and policy corresponding to the supremums.

10 **[Theorem 2]** Step 2 to 3 (line 614) is because $|\mathbb{E}[\cdot]| \leq \mathbb{E}[|\cdot|] \leq \sup |\cdot|$, and step 3 to 4 results from the cancellation
 11 between ω^\top and \arg_Q (as justified above). After line 616, step 2 to 3 arises from the w.l.o.g. assumption that
 12 $\omega^\top Q(s', a', \omega') - \sup_{a'', \omega''} \omega^\top Q'(s', a'', \omega'') \geq 0$, as stated in lines 612 and 615. Thus, the whole expression in $|\cdot|$ is
 13 nonnegative and $\omega^\top Q(s', a', \omega') - \omega^\top Q'(s', a', \omega') \geq 0$. We can discard the last two terms since $\omega^\top Q'(s', a', \omega') \leq$
 14 $\sup_{a'', \omega''} \omega^\top Q'(s', a'', \omega'')$. Step 3 to 4 is because $\sup_{s', \omega'} f(s', a', \omega') \leq \sup_{s', a'', \omega''} f(s', a'', \omega'')$ holds for any a' and $f(\cdot)$.

15 **Empirical results (Reviewers 1 and 3):** **[Multiple runs and error bars]** Each data point in Table 1 indicates the **mean**
 16 **and standard deviation** over **5 independent** training and test runs, for all methods in all four domains. The error bars
 17 in Figure 4 are standard deviations of CR and AE estimated from 5 independent runs under each configuration. This is
 18 mentioned in lines 228, 860-862, 877, 881, but we will consolidate and make this clearer for the reader.

19 **[Statistical tests]** We performed the unpaired t-test between our envelope model and the baselines and achieved
 20 significance scores of $p < 0.05$ vs MOFQI on all domains, $p < 0.01$ vs CN+OLS on DST and $p < 0.05$ vs Scalarized
 21 on FTN, Dialog and SuperMario. We will add this information to the results table.

22 **Comparison with Abels, et al. (Reviewers 2 and 3):** There are 3 key contributions that distinguish our work from
 23 Abels et. al., 2019. We will add a better description of these to the paper as well as better explain figures 2 & 3.

24 **[Algorithmic]** Our algorithm (envelope Q-learning), utilizes the convex envelope of the solution frontier to update
 25 parameters of the policy network, using an optimality filter \mathcal{H} (line 142) which maintains $\sup_{\omega'} \omega^\top Q(\cdot, \cdot, \omega')$. This
 26 allows our method to quickly align one preference with optimal rewards and trajectories that may have been explored
 27 under other preferences. Abels et al. on the other hand, use scalarized updates that optimizes the scalar utility and hence
 28 cannot use the information of $\max_a Q(s, a, \omega')$ to update the optimal solution aligned with a different ω . As illustrated
 29 in Figure 2 (c), assuming we have found two optimal solutions D and F in the CCS, misaligned with preferences ω_2 and
 30 ω_1 . The scalarized update cannot use the information of $\max_a Q(s, a, \omega_1)$ (corresponding to F) to update the optimal
 31 solution aligned with ω_2 or vice versa. It only searches along ω_1 direction leading to non-optimal L, even if solution D
 32 has been seen under ω_2 . Hence, our algorithm has better sample efficiency, as is also seen from the empirical results.

33 **[Theoretical]** Further, we introduce a **theoretical framework** for designing and analyzing value-based MORL al-
 34 gorithms, and **convergence proofs** for our envelope Q-learning algorithm. Abels et al., whose method can also be
 35 analyzed under our framework, do not provide theoretical analyses of the correctness or convergence of their algorithm.

36 **[Empirical]** We also provide **new evaluation metrics and benchmark environments for MORL** – CR and AE. In
 37 terms of experiments, Abels et al. only evaluate on two synthetic domains – DST and Minecart. We apply our algorithm
 38 to a wider variety of domains including DST, FTN and **two complex larger scale domains** – task-oriented dialog and
 39 supermario. Our FTN domain (128 solutions) is a scaled up, more complex version of Minecart (< 10 solutions).