
Empirical Evaluation for Theorem 1 in "Eliciting User Preferences for Personalized Multi-Objective Decision Making through Comparative Feedback"

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This project tests the empirical performance of the algorithm proposed by Shao et al.
2 [2023] in order to find optimal policy for users with unknown preferences through
3 comparative feedback in a Multi-Objective Reinforcement Learning setting. The
4 number of comparative feedback required in practice for user preference estimation
5 meets the theoretical prediction. The optimality gap to the best personalized
6 value confirms the theoretical statement. However, for a simple planning problem,
7 the provided theoretical upper bound for the absolute performance gap is too
8 conservative. The relative performance gap is a better performance evaluation
9 metric. In the worst case scenario, the relative performance optimality gap reaches
10 80%. This poor performance might be caused by the poor preference estimation.
11 Therefore, it is necessary to provide worst case analysis for the relative performance
12 gap and improve user preference estimation.

13 1 Introduction

14 Many real world problems require people to make sequential decisions that balance multiple but
15 sometimes conflicting objectives. For example, in the area of autonomous driving, safety, speed,
16 and comfort are all desired objectives, while speed could negatively impact safety. In classical
17 reinforcement learning, the reward is a scalar that combines several objectives in an arbitrary way.
18 However, different users might have different preferences. Thus, it is important to extend the
19 scalar reward into a vector and design an optimal personalized policy for a given user using as few
20 comparative feedback from them as possible. This paper Shao et al. [2023] provides a provably
21 efficient algorithm to estimate user's personalized policy in a tabular setting. As stated in *Theorem*
22 *1* Shao et al. [2023], given a planning problem, we are able to estimate the personalized optimal
23 value by user's pairwise feedback, with accuracy $O((\sqrt{K} + 1)^{d+\frac{14}{3}} \epsilon^{\frac{1}{3}})$ and $O(K \log(\frac{K}{\epsilon}))$ number
24 of queries, where K is the number of objective, d is the rank of state value matrix, and ϵ is user's
25 comparison distinguishability. However, there is no simulation study in the paper to evaluate the
26 empirical performance of the algorithm. This might cast doubt on the the algorithm's practicality.
27 Therefore, this project aims to bridge the theory-practice gap by adding reinforcement learning
28 experiments to test the algorithm's empirical performance.

29 2 Methodology

30 2.1 Problem Setup

31 If the decision problem has K objectives to consider, then a user is characterized by its preference
 32 vector $\omega^* \in \mathcal{R}^K$, $\|\omega^*\|_1 = 1$. The policy $\pi : |S||A| \rightarrow \Delta_K$ is a mapping from the state action pair to
 33 a discrete distribution vector with size K . The state value function of policy π starting from initial state
 34 s_0 is $V^\pi(s_0) \in \mathcal{R}^K$. $V^\pi = E_{s_0 \sim \rho}[V^\pi(s_0)] \in \mathcal{R}^K$, ρ is the initial state distribution. The personalized
 35 value of policy π is $\langle \omega^*, V^\pi \rangle \in \mathcal{R}^+$. The optimal policy $\pi^* := \arg \max_{\pi \in \Pi} \langle \omega^*, V^\pi \rangle$, and its
 36 corresponding optimal personalized value $\nu^* := \langle \omega^*, V^{\pi^*} \rangle$.

37 2.2 Preference Estimation Algorithm

38 To estimate ω^* , we need three steps, identification of basis policy (see Algorithm 1), computation
 39 of basis ratios (see Algorithm 2), and solving a linear system $\hat{A}\hat{\omega} = e_1$ (see Equation 1). Then the
 estimated optimal policy is $\pi^{\hat{\omega}} = \arg \max_{\pi \in \Pi} \langle \hat{\omega}, V^\pi \rangle$.

```

Initialize  $\pi^{e^*} \leftarrow \pi^{e_1}$ 
for  $j = 2, \dots, k$  do
    compare  $\pi^{e^*}$  and  $\pi^{e_j}$ 
    if  $\pi^{e_j} > \pi^{e^*}$  then  $\pi^{e^*} \leftarrow \pi^{e_j}$ 
end

 $\pi_1 \leftarrow \pi^{e^*}$  and  $u_1 \leftarrow \frac{V^{\pi^{e^*}}}{\|V^{\pi^{e^*}}\|_2}$ 
for  $i = 2, \dots, k$  do
    Arbitrarily pick an orthonormal basis  $\rho_1, \dots, \rho_{k+1-i}$  of  $\text{span}(V^{\pi_1}, \dots, V^{\pi_{i-1}})^\perp$ 
     $j_{\max} \leftarrow \arg \max_{j \in [k+1-i]} \max(|\nu^{\rho_j}|, |\nu^{-\rho_j}|)$ 
    if  $\max(|\nu^{\rho_{j_{\max}}}|, |\nu^{-\rho_{j_{\max}}}|) > 0$  then
         $\pi_i \leftarrow \pi^{\rho_{j_{\max}}}$  if  $|\nu^{\rho_{j_{\max}}}| > |\nu^{-\rho_{j_{\max}}}|$  else  $\pi_i \leftarrow \pi^{-\rho_{j_{\max}}}$ .  $u_i \leftarrow \rho_{j_{\max}}$ 
    end
    else output  $(\pi_1, \pi_2, \dots), (u_1, u_2, \dots)$  and stop
end

```

Algorithm 1: Identification of Basis Policies

40

$$\hat{A} = \begin{bmatrix} V^{\pi_1 T} \\ (\hat{\alpha}_1 V^{\pi_1} - V^{\pi_2})^T \\ \vdots \\ (\hat{\alpha}_{d-1} V^{\pi_1} - V^{\pi_d})^T \end{bmatrix} \quad (1)$$

41 2.3 Multi-Objective Reinforcement Learning (MORL) Environment

42 I extend the Gridworld environment used in HW4 into a MOLR environment. Firstly,
 43 I remove the red-colored self-absorbing states with rewards -10. Secondly, each time,
 44 for user with K objectives, I randomly sample K cells from the grids (corresponds to
 45 $[0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 12, 14, 15, 17, 18, 19, 20, 21, 22, 23, 24]$) except the obstacle states. For
 46 each sampled K cell, I set the reward as 1, to satisfy the bounded reward assumption, $r_t \in [0, 1]$.
 47 I make it as a termination state and adjust the transition matrix T accordingly. I make each cell
 48 correspond to user's objective. The reward function in MORL setting is a $K \times |S| \times |S| \times |A|$ matrix.
 49 In addition, one key assumption in the paper Shao et al. [2023] is the existence of "doing nothing"
 50 policy π_0 and its corresponding value function $V^0 = 0$. To be consistent with this assumption, I add
 51 "STAY" action in addition to the previous four actions, "Move North", "Move South", "Move East",
 52 and "Move West". Therefore, $|A| = 5$. Taking "STAY" action will always receive reward 0. The
 53 noise structure in state transition is the same as the setting in HW4.

Input: $(V^{\pi_1}, \dots, V^{\pi_d})$ and $C_\alpha = 2K$ (see Lemma 1 Shao et al. [2023]), $\hat{\alpha}_i \in [0, C_\alpha], \forall i$
for $i = 1, \dots, d-1$ **do**
 let $l = 0, h = 2C_\alpha, \hat{\alpha}_i = C_\alpha$
 while *True* **do**
 if $\hat{\alpha}_i > 1$ **then**
 compare π_1 and $\frac{1}{\hat{\alpha}_i}\pi_{i+1} + (1 - \frac{1}{\hat{\alpha}_i})\pi_0$. **If** $\pi_1 > \frac{1}{\hat{\alpha}_i}\pi_{i+1} + (1 - \frac{1}{\hat{\alpha}_i})\pi_0$, **then**
 $h \leftarrow \hat{\alpha}_i, \hat{\alpha}_i \leftarrow \frac{l+h}{2}$. **If** $\pi_1 < \frac{1}{\hat{\alpha}_i}\pi_{i+1} + (1 - \frac{1}{\hat{\alpha}_i})\pi_0$, **then** $l \leftarrow \hat{\alpha}_i, \hat{\alpha}_i \leftarrow \frac{l+h}{2}$
 end
 else compare π_{i+1} and $\hat{\alpha}_1\pi_1 + (1 - \hat{\alpha}_i)\pi_0$. **If** $\hat{\alpha}_1\pi_1 + (1 - \hat{\alpha}_i)\pi_0 > \pi_{i+1}$, **then**
 $h \leftarrow \hat{\alpha}_i, \hat{\alpha}_i \leftarrow \frac{l+h}{2}$. **If** $\hat{\alpha}_1\pi_1 + (1 - \hat{\alpha}_i)\pi_0 < \pi_{i+1}$, **then** $l \leftarrow \hat{\alpha}_i, \hat{\alpha}_i \leftarrow \frac{l+h}{2}$.
 if *Indistinguishable* **then**
 break
 end
 end
end
Output $\hat{\alpha}_1, \dots, \hat{\alpha}_{d-1}$

Algorithm 2: Computations of Basis Ratios

54 2.4 Policy Evaluation & Value Iteration

55 Implementing the algorithm requires the knowledge of reinforcement learning. To compare two poli-
56 cies, for example, π_1 and $\frac{1}{\hat{\alpha}_i}\pi_{i+1} + (1 - \frac{1}{\hat{\alpha}_i})\pi_0$, firstly we need to estimate V^{π_1} and $V^{\frac{1}{\hat{\alpha}_i}\pi_{i+1} + (1 - \frac{1}{\hat{\alpha}_i})\pi_0}$.
57 To compute V^π , I use *Policy Evaluation* (see Algorithm 3) and assume initial state following a uniform
58 distribution. Secondly, based on the estimated state value function, user returns "indistinguishable "
59 if $|\langle \omega^*, V^{\pi_1} \rangle - \langle \omega^*, V^{\frac{1}{\hat{\alpha}_i}\pi_{i+1} + (1 - \frac{1}{\hat{\alpha}_i})\pi_0} \rangle| < \epsilon$. Otherwise, user will choose the policy with
60 higher personalized value. π_1, \dots, π_d are greedy policies with respect to $V^{\pi_1}, \dots, V^{\pi_d}$ accordingly.
61 They are jointly estimated by *Value Iteration* (see Algorithm 4).

Input: $\theta > 0$ tolerance parameter, γ discount factor, π , policy to evaluation

Initialize $V(s)$ arbitrarily, with $V(\text{terminal}) = 0$

Repeat:

$\Delta \leftarrow 0$

for $s \in S$ **do**

62 $v \leftarrow V(s)$
 $V(s) \leftarrow \sum_{a \in A(s)} \pi(a|s) \sum_{s', r} p(s', r|s, a)(r + \gamma V(s'))$
 $\Delta \leftarrow \max(\Delta, |v - V(s)|)$

end

Until $\Delta < \theta$

Output: $V^\pi(s), \forall s$

Algorithm 3: Policy Evaluation

Input: $\theta > 0$ tolerance parameter, γ discount factor

Initialize $V(s)$ arbitrarily, with $V(\text{terminal}) = 0$

Repeat:

$\Delta \leftarrow 0$

for $s \in S$ **do**

63 $v \leftarrow V(s)$
 $V(s) \leftarrow \max_{a \in A(s)} \sum_{s', r} p(s', r|s, a)(r + \gamma V(s'))$
 $\Delta \leftarrow \max(\Delta, |v - V(s)|)$

end

Until $\Delta < \theta$

Output: $V(s), \forall s$ and π : greedy policy w.r.t. $V(s)$

Algorithm 4: Value Iteration

Parameter	θ	γ	noise	ϵ	K	rep
Description	tolerance threshold	discounting factor	-	indistinguishability	no.of objectives	-
Value	0.01	0.99	0.1	0.01	[3, 10]	5

Table 1: Experiment Parameters

2.4.1 Experiment Setup

3 Result & Discussion

3.1 Sanity Check

In Figure 1, we set $K = 3$. The user has three objectives, going to the green cell, going to the blue cell, and going to the yellow cell. The left column shows state value estimates ($V^{\pi^T(s)}_{\omega}, \forall s \in S$) and corresponding optimal actions in the MORL setting when we input the preference vector $[1, 0, 0], [0, 1, 0], [0, 0, 1]$. This corresponds to prioritizing going to the green cell, the blue cell, and the yellow cell. The right column shows state value estimates and associated optimal actions in the GridWorld setting with single objective and scalar reward. We can see two settings having the same estimates and optimal actions. This confirms the correctness of *Value Iteration* implementation in the MORL setting. The last subplot in Figure 1 shows the state value estimates of "doing nothing" policy. In the implementation, $\pi_0 \in \mathcal{R}^{|S| \times |A|}, \pi_0[:, \text{STAY}] = 1$. This policy is supposed to have $V_0 = 0$ and this is what we get. This confirms the correctness of *Policy Evaluation* implementation in the MORL setting.

3.2 Performance Evaluation

We want to evaluate the empirical performance of Theorem 1. The full statement Shao et al. [2023] is

Theorem 1 Consider the algorithm of computation \hat{A} defined in Eq(1) and any solution $\hat{\omega}$ to $\hat{A}x = e_1$ and outputting the policy $\pi^{\hat{\omega}} = \arg \max_{\pi \in \Pi} \langle \hat{\omega}, V^{\pi} \rangle$, which is the optimal personalized policy for preference vector $\hat{\omega}$. Then the output policy $\pi^{\hat{\omega}}$ satisfying that $\nu^* - \langle \omega^*, V^{\pi^{\hat{\omega}}} \rangle \leq O((\sqrt{K} + 1)^{d + \frac{14}{3}} \epsilon^{\frac{1}{3}})$ by using $O(K \log(\frac{K}{\epsilon}))$ comparison queries.

In Figure 2, the red dotted line plots the theoretical upper bound $O(K \log(\frac{K}{\epsilon}))$ of the number of comparison queries required to estimate ω^* . The black dotted line shows the average number of comparison queries required in 5 repeated trials for each $K \in [3, 10]$ simulated under different seeds. The shaded blue region indicates the minimum and maximum number of comparison queries for each K . None of the required number of comparison queries exceeds the theoretical upper bound. Practice meets theory. Since the maximum number of required queries hits the upper bound when $K = 5$, this suggests that the theoretical upper bound is tight.

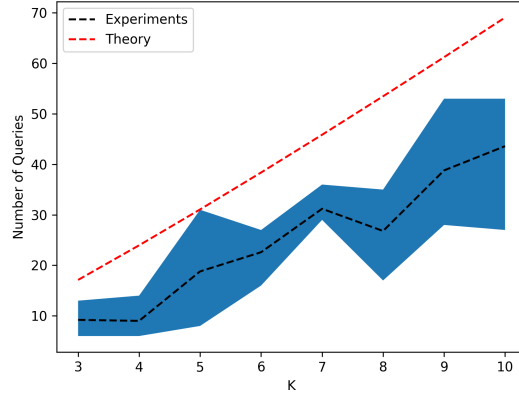


Figure 2: Simulation Results of Number of Comparison Queries Required in Practice to Estimate ω^*

Figure 3 shows the empirical performance of the absolute performance gap defined as $\nu^* - \langle \omega^*, V^{\pi^{\hat{\omega}}} \rangle$. Similar to Figure 2, for each K , the black dotted line shows the average performance over 5 trials and the shaded blue region denotes the minimum and maximum. When $K = 3, d = 1$, the theory (see Theorem 1) predicts the upper bound as $(\sqrt{K} + 1)^{d + \frac{1}{3}} \epsilon^{\frac{1}{3}} = 64 \geq 1$. It is increasing in terms of K and d . Therefore, the absolute performance gap in practice automatically meets the theoretical statement for all K . However, we do not observe an obvious increase in the absolute performance gap in terms of K and d , as opposed to the theoretical prediction.

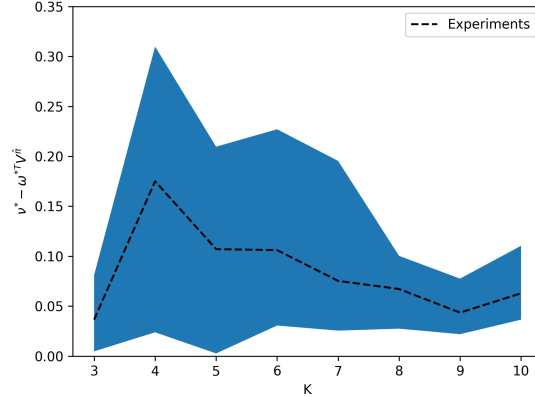


Figure 3: Simulation Results of Absolute Performance Gap in Practice

Since the implemented MORL is very simple, the theoretical upper bound for the absolute performance gap is too loose. In this case, the relative performance gap, $\frac{\nu^* - \langle \omega^*, V^{\pi^{\hat{\omega}}} \rangle}{\nu^*}$ should be a better metric to evaluate the actual performance of the proposed algorithm. As Figure 4 suggests, even though in the best case, the optimality gap is around 10% to 20%, in the worst case, however, the relative performance gap could be as huge as 80%. This signifies the importance of understanding and bounding the relative performance gap in the worst case scenario, especially when we have simple planning problems.

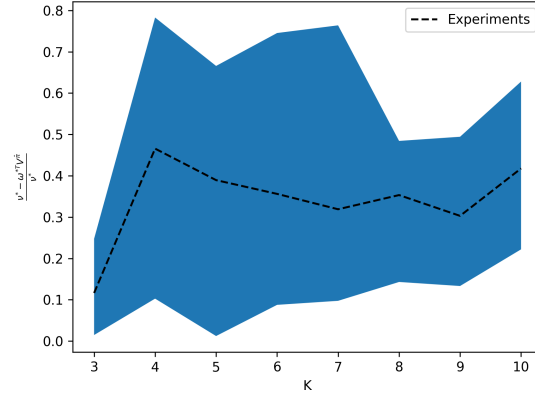


Figure 4: Simulation Results of Relative Performance Gap in Practice

3.3 Preference estimation

To understand which factor is associated with the relative performance gap, I plot the randomly generated true preference and the estimated preference for $K = 3$ (see Figure 5). Due to the space limit and human interpretability constraint, I omit showing plots for $K > 3$. These plots can be reproduced here. For trial number 0, the true preference is $[0.33, 0.33, 0.34]$, while the estimated preference is $[0, 0, 1]$. These two preference vectors look very different but they have very similar personalized values. This is because when $\omega > 0$, the three objectives going to green, blue, or yellow cells, are exchangeable. Situations get more complicated, however, when ω includes negative values, i.e., users prefer going to the yellow cell while avoiding going to the blue cell. Intuitively, we expect good sign alignment in each coordinate of ω^* and $\hat{\omega}$. Unfortunately, Table 2 shows that this conjecture is not always true. The worst case could happen when only 1 out of 5 signs aligns between the true preference and the estimated preference. Even though the personalized values between the true preference and the estimated preference could be similar, it is counter intuitive for a preference estimation with many "prefer / hate" flips to be a good estimation and has low optimality gap.

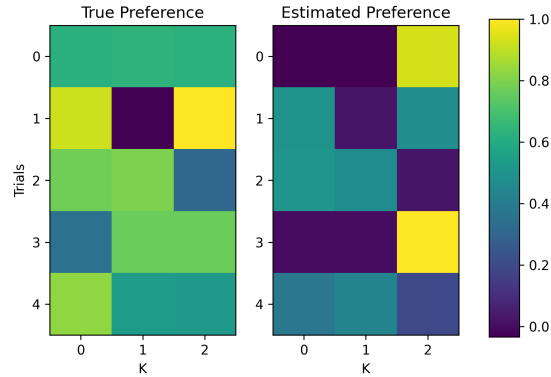


Figure 5: Compare True Preference and Estimated Preference

K	3	4	5	6	7	8	9	10
Mean	2.2	3.8	3.8	5.2	4.6	7.2	6.6	7.4
Min	1	3	1	4	3	5	4	7
Max	3	4	5	6	7	8	9	9

Table 2: Counts of Aligned Sign between True Preference & Estimated Preference

119 4 Conclusion

120 In conclusion, empirical assessment does not show any violation of *Theorem 1*. For the implemented
121 simple MORL, the theoretical bound of the absolute performance gap is too loose. Instead, the
122 relative performance gap is a better performance evaluation metric. In the worst case scenario, the
123 relative performance optimality gap could be as large as 80%. This poor performance might be
124 caused by the poor preference estimation assessed by the proportion of flipped signs. Based on these
125 observations, I suggest Shao et al. [2023] to consider two additional analysis. Firstly, provides worst
126 case analysis for the relative performance gap. Secondly, bounds the difference between $\hat{\omega}$ and ω^* by
127 putting necessary restrictions on $\hat{\omega}$. If $d \ll K$, the solution space of $\hat{A}x = e_1$ might be very too
128 large, which makes $\hat{\omega}$ too flexible to be true.

129 5 Limitations

130 This project has several limitations. First, the implemented MORL environment is almost the
131 simplest tabular MDP. Given its simple structure, it might not be able to test the robustness of the
132 $O((\sqrt{K} + 1)^{d + \frac{14}{3}} \epsilon^{\frac{1}{3}})$ in the worse case, especially when H is large. In addition, randomly sampling
133 K cells from the grid and setting their rewards as 1 is not a good way to model K different objectives.
134 For the implemented environment, K objectives are exchangeable when $\omega > 0$. In reality, however,
135 K objectives can conflict with each other and in general, they are not exchangeable.

136 References

137 H. Shao, L. Cohen, A. Blum, Y. Mansour, A. Saha, and M. R. Walter. Eliciting user preferences for
138 personalized multi-objective decision making through comparative feedback, 2023.

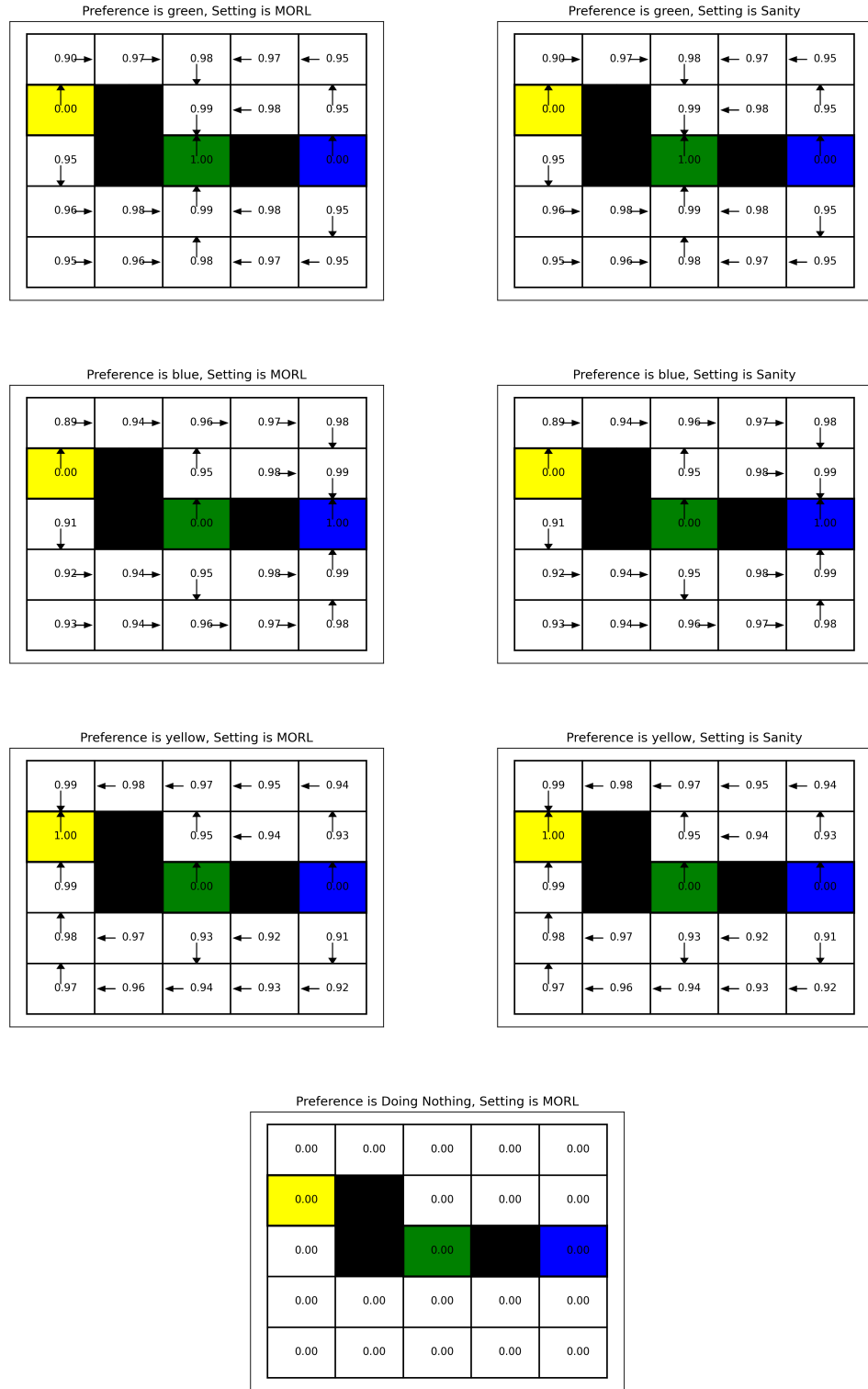


Figure 1: Sanity Check