

1 We would like to thank the reviewers for appreciating our novel contributions on the algorithmic and theoretical front!
2 We focus on clarifying our experimental results in this rebuttal.

3 [*Why DM fails at ModelFail and SSD-IS achieve EXACTLY the same results as DM at ModelFail?*].

4 ModelFail was first introduced by Thomas and Brunskill [2016] to show the failure of model-based approach in the
5 MDPs with some partial observability. In ModelFail, the agent cannot tell the difference between any of the states
6 except for s_1 , but both DM and SSD-IS require full observability. From the point of view of both DM and SSD-IS, the
7 actions have no impact on state transitions or rewards, so every policy has the same cumulative reward (equal to the
8 true cumulative reward of the behavior policy). A detailed discussion about why DM fails at ModelFail can be found in
9 [Thomas and Brunskill, 2016, Section D.1]. MIS can handle partial observability by using observable states and the
10 partial trajectories between them. Please refer Section 5.1 (line 258-262, there is a typo in Line 262, $\frac{\pi(a_{2\tau}^{(i)}|s_{2\tau}^{(i)})}{\mu(a_{2\tau}^{(i)}|s_{2\tau}^{(i)})}$ should
11 be $\frac{\pi(a_{2\tau}^{(i)}|?)}{\mu(a_{2\tau}^{(i)}|?)}$, where symbol “?” stands for “unobserved”, is an observed variable that the policy needs to react upon).
12 Also see Section C (line 567-575) in the supplement for more details.

13 [*Why MIS outperforms SSD-IS in time-invariant environments (including MountainCar) when n is large?*].

14 The time-invariant ModelWin and MountainCar we used in the paper are finite-horizon *undiscounted* MDPs. Even
15 though these environments have time-invariant transitions, the state marginal distributions at each t actually change
16 with time and only converge to the stationary distribution as $t \rightarrow \infty$.

17 SSD-IS uses the stationary distribution ($t \rightarrow \infty$) to approximate that for all $t = 1, \dots, H$ which is biased and not
18 consistent even as the number of episodes $n \rightarrow \infty$. MIS, on the other hand, uses nearly unbiased and consistent
19 estimators of the state marginals at every t . This allows MIS to outperform SSD-IS on Mountain Car when n gets large.
20 We believe this is the reason and we will investigate it in details in our future work.

21 **Reviewer #1**

22 [*“A specific baseline I would really like to see the authors add is the PDIS (per-decision IS) and CWPDIS (consistent
23 weighted per-decision IS).”*]

24 The IS and WIS in the experiments are step-wise, which are essentially PDIS and CWPDIS. The detailed explanation is
25 in Section 3 and Section C.

26 [*“Why does it (SSD-IS) achieve ... perform as well as MIS for mountain car but eventually stops improving?”*]

27 Please check the answers at the beginning.

28 [*“If p_t is sampled uniformly at each time step, isn’t ... setting equivalent to a time-invariant MDP with $p = 3.5$?”*]

29 Sorry for the confusion. Note that each transition probability p_t is only sampled before the experiments and fixed during
30 the experiments for all episodes. We will clarify it in the final version.

31 **Reviewer #2**

32 Thanks for supporting our paper. We are planning to extend our approach to large-scale environments with extensive
33 function approximation.

34 **Reviewer #3**

35 [*“In Figure 2 and 3, why DM and SSD-IS method works well in ModelWin but perform very bad at ModelFail?”*]

36 [*“For me it is surprised in time-invariant environment SSD-IS method perform worse than MIS method.”*]

37 Please check the answers at the beginning.

38 [*“In Figure 3 (b) and (d), why the curve is not smooth even after 128 repetition?”*]

39 Note that the Y-axis is relative MSE, which is normalized by the true cumulative reward. In this time-varying MDP
40 (Figure 3), the true cumulative reward is related to the transition probabilities p_t at each time step. We sample each
41 p_t before the experiments and then fix them during the experiments, so the true cumulative reward is a non-smooth
42 function of H and the figures with increasing H should not be smooth. In the time-invariant MDP (Figure 2), the true
43 cumulative reward is a smooth function of H and the corresponding figures are smooth.

44 **References**

45 Thomas, P. and Brunskill, E. (2016). Data-efficient off-policy policy evaluation for reinforcement learning. In
46 *International Conference on Machine Learning*, pages 2139–2148.