

1 Our paper presents an geometric interpretation for inverse reinforcement learning (IRL) with finite states and actions,
 2 as well as a corresponding L1-regularized Support Vector Machine formulation with formal guarantees in sample
 3 complexity and with regard to Bellman optimality. We thank the reviewers for their feedback and for bringing several
 4 issues and improvements to our attention. Our main contribution and focus in this paper is theoretical, which the
 5 reviewers seem to agree with, and the experiments serve primarily to validate our theory.

6 Following comments from the reviewers, we will move The proof of Theorem 5.1 to the appendix and add content
 7 discussing the background of IRL. In addition to the linear programming method [5] and the Bayesian IRL method [6]
 8 presented in the original submission, several other approaches to the IRL problem exist. These include Hybrid IRL
 9 [4], Maximum Margin Planning [7], Maximum Entropy IRL (MaxEnt) [9], and Gaussian Process IRL (GPIRL) [3]. A
 10 survey of other approaches to the inverse reinforcement learning problem can be found in [1]. While the problem in our
 11 paper is formulated in the form of a standard Markov Decision Process (MDP), several approaches instead consider
 12 the linearly-solvable MDP (LMDP) formulation presented in [2]. As mentioned in [2] (in particular Section 2.6), the
 13 problem formulation of the standard MDP and LMDP is different. Given the true transition probabilities, methods such
 14 as Multiplicative Weights for Apprenticeship Learning (MWAL) [8] and [5] are guaranteed to recover the true optimal
 15 policy where as methods that use LMDP to solve the same problem are not. To confirm this, we have used GPIRL [3]
 16 code available at <https://graphics.stanford.edu/projects/gpir1/> to solve the same synthetic experiments
 17 presented in Figure 3 of Section 8 of our paper and found the rewards from GPIRL were not Bellman optimal as per the
 18 definition in section 3 line 75 of our paper. Further comparisons with GPIRL, Bayesian IRL and MWAL are shown
 19 the Figure 1. We note that about 30% of the rewards returned by MWAL were the trivial solution $R = 0$ which were
 20 not counted as successes. The result presented in our paper immediately impacts algorithms that use standard MDP
 21 models more it impacts than algorithms that use LMDP models such as MaxEnt and GPIRL, as the objective of the
 22 LMDP-based algorithms is different. In the case of standard MDP problems it readily provides a sample complexity
 23 and a formal guarantee with respect to Bellman optimality, which is not provided by any of the other methods.

24 The formulation provided in our paper is a nonparametric approach as compared to approaches that use features derived
 25 from states. It also makes no assumptions on the sparseness of the transitions. Our experiments reflect this as well,
 26 as the transition probabilities are drawn from a uniform distribution with no sparseness assumptions and would be
 27 more difficult to than sparse cases. In contrast, tasks like gridworld tend to have have sparse transitions probabilities
 28 or feature transformations that reduce dimensions. We provide a method and a guarantee on optimality that does not
 29 require sparseness and does not depend on feature selection, which is a problem with other methods.

30 Definition 4.1 leads to our paper considering only Regime 3 cases. From a practical perspective this is not
 31 a loss. All problems within Regime 1 have only one feasible solution, $R \equiv 0$ which does not provide
 32 any information with respect to Bellman optimality as no policy is preferred over the other with this reward.
 33 Problems in Regime 2 require an infinite number of samples to both ascertain that
 34 the problem is indeed in Regime 2 and to solve the problem as perturbations can
 35 lead to the estimated problem falling under Regime 1 or Regime 3 as mentioned in
 36 line 134 of our paper. As stated in line 130, this type of assumption has been made
 37 in other areas as well.

38 While our paper does state the parameters and objective of the finite state MDP IRL
 39 problem from lines 39-76 as well as an initial formulation, our paper does not aim
 40 to derive the results and formulation of [5] from scratch. As suggested, we will
 41 present this problem as a standalone problem and explain the origin of the F matrix
 42 from Bellman Optimality.

43 References

44 [1] S. Arora and P. Doshi. A survey of inverse reinforcement learning: Challenges, methods and progress. *arXiv:1806.06877*, 2018.

45 [2] K. Dvijotham and E. Todorov. Inverse optimal control with linearly-solvable mdps. In *ICML 2010*, pages 335–342, 2010.

46 [3] Sergey Levine, Zoran Popovic, and Vladlen Koltun. Nonlinear inverse reinforcement learning with gaussian processes. In *NeurIPS*
 47 *2011*, pages 19–27, 2011.

48 [4] Gergely Neu and Csaba Szepesvári. Apprenticeship learning using inverse reinforcement learning and gradient methods. In *UAI*
 49 *2007*, pages 295–302. AUAI Press, 2007.

50 [5] A. Y. Ng and S.J. Russel. Algorithms for inverse reinforcement learning. In *ICML 2000*, pages 663 – 670, 2000.

51 [6] Deepak Ramachandran and Eyal Amir. Bayesian inverse reinforcement learning. *IJCAI 2007*, 51(61801):1–4, 2007.

52 [7] N. D. Ratliff, J. A. Bagnell, and M. A. Zinkevich. Maximum margin planning. In *ICML 2006*, pages 729–736. ACM, 2006.

53 [8] U. Syed, M. Bowling, and R. Schapire. Apprenticeship learning using linear programming. In *ICML 2008*, pages 1032–1039. ACM,
 54 2008.

55 [9] Brian D Ziebart, Andrew Maas, J Andrew Bagnell, and Anind K Dey. Maximum entropy inverse reinforcement learning. *AAAI*
 56 *2008*.

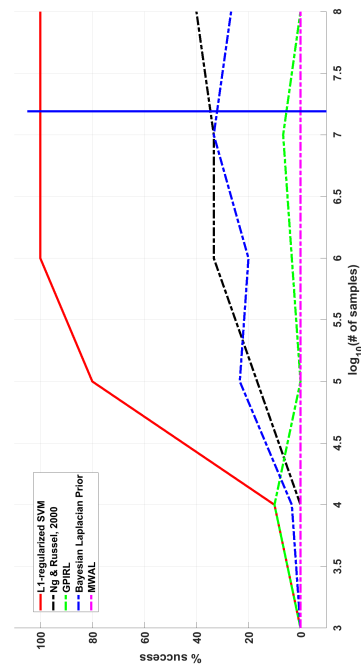


Figure 1: