
Reward-Free Reinforcement Learning with GNN and Adversarial Linear Mixture MDPs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Reward-free RL is independently developed in the unconstrained literature, which
2 learns the transition dynamics without using the reward information and is thus
3 naturally capable of addressing RL with multiple objectives under the common
4 dynamics. This paper proposes a new framework for the reward-free RL setting
5 with function approximation i.e. the adversarial linear mixture MDPs. As Jin, et al.
6 (2020). We partition this setting into an exploration phase and a planning phase.
7 During the exploration phase, the agent first collects trajectories from an MDP
8 M without a pre-specified reward function. Using the Graph Neural Networks
9 (GNNs) to store the significant states in dataset D instead of all states, each with a
10 heuristic weight. In the planning phase, it is tasked with computing near-optimal
11 policies under M for a collection of given reward functions. The agent generalizes
12 previously learned information using the linear mixture MDPs that allows it to
13 approximate the policy given an arbitrary reward function.

14 1 Introduction

15 In reinforcement learning (RL), an agent repeatedly interacts with an unknown environment with the
16 goal of maximizing its cumulative reward. To do so, the agent must engage in exploration, learning
17 to visit states to investigate whether they hold high rewards.

18 Exploration is widely regarded as the most significant challenge in RL, because the agent may have
19 to take precise sequences of actions to reach states with high reward. Here, simple randomized
20 exploration strategies provably fail: for example, a random walk can take exponential time to reach
21 the corner of the environment where the agent can accumulate high rewards (Li, 2012). While
22 reinforcement learning has seen a tremendous surge of recent research activity, essentially all of
23 the standard algorithms deployed in practice employ simple randomization or its variants, and
24 consequently incur extremely high sample complexity.

25 In this extended abstract paper, we aim to develop an end-to-end instantiation of this proposal. To
26 this end we ask: How can we generalize the concepts of significant states and coverage guarantees?
27 And how can we develop such an agent that can generalize enough?

28 1.1 Notations

29 In the reward-free setting, we would like to design algorithms that efficiently explore the state space
30 without the guidance of reward information. Over the course of K episodes, the agent collects a
31 dataset of visited states, actions, and transitions $D = s_h^{(k)}, a_h^{(k)}(k, h) \in [k] \times [H]$, which is the
32 outcome of the exploration phase.

33 **Graph Neural Networks** Graph neural networks (GNN) are a class of neural networks that operate
34 directly on graph-structured data. A wide variety of graph neural network architectures have been

35 proposed. These range from simple graphs, to directed graphs, to graphs that contain information,
 36 up to convolutional graphs. The graph $G = (N, E)$ is defined as having nodes $n_i \in N$ and directed
 37 edges $e_{ij} \in E$ from node n_j to n_i . Both – the nodes and the edges – contain additional information.
 38 The node value is denoted as h_i for the i -th node and the edge value as e_{ij} connecting the i -th with
 39 the j -th node. In each layer of the GNN, a dense node neural network layer is applied per node and a
 40 dense edge neural network layer per edge. Each GNN layer has three computation steps: First, the
 41 next edge values e_{ij}^{k+1} are computed using the current edge values e_{ij}^k , the from-node values h_i^k and
 42 the to-node values h_j^k . These values are concatenated and passed into a dense neural network layer
 43 $f_x^k(\cdot)$ that is parameterized by X . This can be represented as:

$$e_{ij}^{k+1} = f_x^k([h_i^k, e_{ij}^k, h_j^k]) \quad (1)$$

44 **Linear Mixture MDPs.** We focus on a special class of MDPs named linear mixture MDPs (Ayoub
 45 et al., 2020; Cai et al., 2020; Zhou et al., 2021; He et al., 2022; Li et al., 2023), where the transition
 46 kernel is linear in a known feature mapping $\phi : SAS \rightarrow R^d$ with the following definition.

47 Definition 1 (Linear Mixture MDPs). An MDP instance $M = (S, A, H, P_{hh=1}^H), l_{kk=1}^K$ is called
 48 an inhomogeneous, episodic B-bounded linear mixture MDP if there exists a known feature mapping
 49 $\phi(s'|s, a) : SAS \rightarrow R^d$ with $\phi(s'|s, a) \geq 1$ and unknown vectors $\phi_{hh=1}^* \in R^d$ with $\phi_h^* \geq B$ such that
 50 for all $(s, a, s') \in SAS$ and $h \in [H]$, it holds that $P_h(s'|s, a) = \langle \phi(s'|s, a), \phi_h^* \rangle$

51 2 Approximate MDP Solvers

52 Approximate MDP solvers aim to find a near-optimal policy when the exact transition matrix P
 53 and reward r are known. The simplest way to achieve this is by the Value Iteration (VI) algorithm,
 54 which solves the Bellman optimality equation in a dynamical programming fashion. Then the greedy
 55 policy induced by the result Q^* gives precisely the optimal policy without error. Another popular
 56 approach frequently used in practice is the Natural Policy Gradient (NPG) algorithm. In each iteration,
 57 the algorithm first evaluates the value of policy $\pi^{(t)}$ using Bellman equation. Then it updates the
 58 policy by first scaling it with the exponential of learning times value $Q^{\pi^{(t)}}$, and then performs a
 59 normalization. For completeness, we provide its guarantee here, which resembles the infinite horizon
 60 analysis in (Agarwal et al., 2019)

61 References

- 62 [1] Chi Jin, Akshay Krishnamurthy, Max Simchowitz, and Tiancheng Yu. Reward-free exploration for reinforce-
 63 ment learning. In International Conference on Machine Learning, pages 4870–4879. PMLR, 2020.
- 64 [2] Li, L. Sample complexity bounds of exploration. In Reinforcement Learning, pp. 175–204. Springer, 2012
- 65 [3] Ayoub, A., Jia, Z., Szepesv ´ari, C., Wang, M., and Yang, L. (2020). Model-based reinforcement learning with
 66 value-targeted regression. In Proceedings of the 37th International Conference on Machine Learning (ICML),
 67 pages 463–474.
- 68 [4] Cai, Q., Yang, Z., Jin, C., and Wang, Z. (2020). Provably efficient exploration in policy optimization. In
 69 Proceedings of the 37th International Conference on Machine Learning (ICML), pages 1283–1294.
- 70 [5] Zhou, D., Gu, Q., and Szepesv ´ari, C. (2021). Nearly minimax optimal reinforcement learning for linear
 71 mixture Markov decision processes. In Proceedings of the 34th Conference on Learning Theory (COLT), pages
 72 4532–4576.
- 73 [6] He, J., Zhou, D., and Gu, Q. (2022). Near-optimal policy optimization algorithms for learning adversarial
 74 linear mixture MDPs. In Proceedings of the 25th International Conference on Artificial Intelligence and Statistics
 75 (AISTATS), pages 4259–4280.
- 76 [7] Li, L.-F., Zhao, P., and Zhou, Z.-H. (2023). Dynamic regret of adversarial linear mixture MDPs. In Advances
 77 in Neural Information Processing Systems 36 (NeurIPS), pages 60685–60711.
- 78 [8] Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. Optimality and approximation with policy gradient
 79 methods in markov decision processes. arXiv preprint arXiv:1908.00261, 2019.