

---

# Active Transfer Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1

## 2 Introduction

## 3 Preliminaries

4 **Markov Decision Processes** We define a Markov decision process (MDP) as a tuple  $\mathcal{M} =$   
5  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, p_0, \gamma \rangle$ , where  $\mathcal{S}$  is the state-space,  $\mathcal{A}$  is a finite set of actions,  $\mathcal{P}(\cdot|s, a)$  is the distribution  
6 of the next state  $s'$  given that action  $a$  is taken in state  $s$ ,  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function,  $p_0$   
7 is the initial-state distribution, and  $\gamma \in [0, 1)$  is the discount factor. We assume the reward function  
8 to be uniformly bounded by a constant  $R_{max} > 0$ . A deterministic policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  is a mapping  
9 from states to actions. At the beginning of each episode of interaction, the initial state  $s_0$  is drawn  
10 from  $p_0$ . Then, the agent takes the action  $a_0 = \pi(s_0)$ , receives a reward  $\mathcal{R}(s_0, a_0)$ , transitions to the  
11 next state  $s_1 \sim \mathcal{P}(\cdot|s_0, a_0)$ , and the process is repeated. The goal is to find the policy maximizing  
12 the long-term return over a possibly infinite horizon:  $\max_{\pi} J(\pi) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t | \mathcal{M}, \pi]$ . To this  
13 end, we define the optimal value function  $Q^*(s, a)$  as the expected return obtained by taking action  
14  $a$  in state  $s$  and following an optimal policy thereafter. Then, an optimal policy  $\pi^*$  is a policy that  
15 is greedy with respect to the optimal value function, i.e.,  $\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$  for all states  
16  $s$ . It can be shown (e.g., [4]) that  $Q^*$  is the unique fixed-point of the optimal Bellman operator  $T$   
17 defined by  $TQ(s, a) = \mathcal{R}(s, a) + \gamma \mathbb{E}_{\mathcal{P}}[\max_{a'} Q(s', a')]$  for any value function  $Q$ . From now on, we  
18 adopt the term  $Q$ -function to denote any plausible value function, i.e., any function  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$   
19 uniformly bounded by  $\frac{R_{max}}{1-\gamma}$ .

20 We define the Bellman error (or Bellman residual) of a  $Q$ -function  $Q$  as  $B(Q) \triangleq TQ - Q$ . Notice  
21 that a  $Q$ -function  $Q$  is optimal if, and only if,  $\|B(Q)\|_{\infty} = 0$ .

22 **Multitask Settings** We represent tasks  $\tau$  as MDPs with shared state and action spaces, but with  
23 potentially different values for all other parameters. We assume the existence of a distribution  $\mathcal{D}$  over  
24 tasks, i.e.,  $\tau \sim \mathcal{D}$ , and we suppose that we are able to sample from such distribution.

## 25 3 Approach

26 We start by noticing that the task distribution  $\mathcal{D}$  clearly induces a distribution over optimal  $Q$ -  
27 functions. Then, our goal is to estimate such distribution from the set of source tasks and use it as a  
28 prior for speeding up the learning process in the target task.

29 We assume states and actions to be drawn from a fixed joint distribution  $\mu$ . We define the set of  
30  $Q$ -functions of our interest as the set  $\mathcal{Q}^{\epsilon}$  of all  $Q$  functions whose Bellman error  $B(Q)$  is in some  
31  $\epsilon$ -ball defined by the  $l_p$ -norm  $\|\cdot\|_{p, \mu}$ :

$$\mathcal{Q}^{\epsilon} = \{Q \in \mathcal{Q} \mid \|B(Q)\|_{p, \mu}^p \leq \epsilon\} \quad (1)$$

32 Given a dataset of  $N$  samples  $D = \langle s_i, a_i, r_i, s'_i \rangle_{i=1}^N$ , we can approximate the  $l_p$ -norm of the Bellman  
 33 error of a  $Q$ -function  $Q$  as:

$$\|B(Q)\|_{p,D}^p = \frac{1}{N} \sum_{i=1}^N \left| r_i + \gamma \max_{a'} Q(s'_i, a') - Q(s_i, a_i) \right|^p \quad (2)$$

34 **Theorem 1.** Let  $Q$  be a  $Q$ -function with empirical Bellman error, computed on a dataset  $D$  of  $N$   
 35 i.i.d. samples, given by  $\|B(Q)\|_{p,D}^p = \hat{q}$ . Then, for any  $\epsilon \geq 0$ :

$$P\left(Q \in \mathcal{Q}^\epsilon \mid \|B(Q)\|_{p,D}^p = \hat{q}\right) \leq \exp\left(-\frac{2N \max\{\epsilon - \hat{q}, 0\}^2}{\left(\frac{2R_{max}}{1-\gamma}\right)^{2p}}\right) \quad (3)$$

36 *Proof.* Assume  $\hat{q} > \epsilon$ . Then:

$$\begin{aligned} P\left(Q \in \mathcal{Q}^\epsilon \mid \|B(Q)\|_{p,D}^p = \hat{q}\right) &= P\left(\|B(Q)\|_{p,\mu}^p \leq \epsilon \mid \|B(Q)\|_{p,D}^p = \hat{q}\right) \\ &= P\left(\|B(Q)\|_{p,\mu}^p - \hat{q} \leq \epsilon - \hat{q} \mid \|B(Q)\|_{p,D}^p = \hat{q}\right) \end{aligned}$$

37 Notice that  $\mathbb{E}[\hat{q}] = \|B(Q)\|_{p,\mu}^p$  and that  $\epsilon - \hat{q} < 0$  by assumption. Then, we can apply Hoeffding's  
 38 inequality to write:

$$P\left(Q \in \mathcal{Q}^\epsilon \mid \|B(Q)\|_{p,D}^p = \hat{q}\right) \leq \exp\left(-\frac{2N(\epsilon - \hat{q})^2}{\left(\frac{2R_{max}}{1-\gamma}\right)^{2p}}\right)$$

39 Finally, when  $\hat{q} \leq \epsilon$ , the probability can be straightforwardly upper-bounded by 1. Combining the  
 40 two results concludes the proof.  $\square$

41 **Theorem 2.** Let  $Q$  be a  $Q$ -function with empirical Bellman error, computed on a dataset  $D$  of  $N$   
 42 i.i.d. samples, given by  $\|B(Q)\|_{p,D}^p = \hat{q}$ . Then, for any  $\epsilon \geq 0$ :

$$P\left(\|B(Q)\|_{p,D}^p = \hat{q} \mid Q \in \mathcal{Q}^\epsilon\right) \leq \frac{\epsilon}{\hat{q}}$$

*Proof.*

$$\begin{aligned} P\left(\|B(Q)\|_{p,D}^p = \hat{q} \mid Q \in \mathcal{Q}^\epsilon\right) &= P\left(\|B(Q)\|_{p,D}^p = \hat{q} \mid \|B(Q)\|_{p,\mu}^p \leq \epsilon\right) \\ &\leq P\left(\|B(Q)\|_{p,D}^p \geq \hat{q} \mid \|B(Q)\|_{p,\mu}^p \leq \epsilon\right) \\ &\leq \frac{E[\|B(Q)\|_{p,D}^p]}{\hat{q}} \\ &\leq \frac{\epsilon}{\hat{q}} \end{aligned}$$

43 The first inequality is straightforward, the second one is from Markov's inequality, while the third  
 44 one is due to the fact that  $Q \in \mathcal{Q}^\epsilon$ .  $\square$

## 45 4 Related Works

46 List of any paper that relates to our approach together with a brief description:

- 47 • [3]: the authors propose a method for efficient exploration via randomized value functions.  
 48 The optimal  $Q$ -function is computed by bayesian LSVI and, after each update, parameters  
 49 are sampled from the posterior. Then, the agent follows a greedy policy with respect to the  
 50 sampled  $Q$ -function. Regret bounds are provided.

- 51 • [2]: the authors extend the idea of randomized value functions to drive exploration in deep  
52 RL. A posterior distribution over  $Q$ -functions is approximated via bootstrapping. In each  
53 episode, the agent acts greedily with respect to a  $Q$ -function sampled from the approximated  
54 posterior.
- 55 • [1]: the authors build on top of bootstrapped DQNs to provide UCB-like exploration bonuses.  
56 In a previous (?) version of the paper, exploration bonuses based on information gain are  
57 also proposed.

## 58 **5 Experiments**

## 59 **6 Conclusion**

## 60 **References**

- 61 [1] Richard Y Chen, John Schulman, Pieter Abbeel, and Szymon Sidor. Ucb and infogain exploration via  
62 q-ensembles. *arXiv preprint arXiv:1706.01502*, 2017.
- 63 [2] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped  
64 dqn. In *Advances in neural information processing systems*, pages 4026–4034, 2016.
- 65 [3] Ian Osband, Benjamin Van Roy, and Zheng Wen. Generalization and exploration via randomized value  
66 functions. *arXiv preprint arXiv:1402.0635*, 2014.
- 67 [4] Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley  
68 & Sons, Inc., New York, NY, USA, 1994.

