

Target Transfer Q-Learning and Its Convergence Analysis

Yue Wang^{†*}, Qi Meng[‡], Wei Cheng[‡], Yuting Liug[†], Zhi-Ming Ma^{†§}, Tie-Yan Liu[‡]

[†]School of Science, Beijing Jiaotong University, Beijing, China {11271012, ytliu}@bjtu.edu.cn

[‡]Microsoft Research, Beijing, China {meq, wche, Tie-Yan.Liu}@microsoft.com

[§] Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing, China mazm@amt.ac.cn

Abstract

Reinforcement Learning (RL) technologies are powerful to learn how to interact with environments and have been successfully applied to variants of important applications. Q-learning is one of the most popular methods in RL, which uses temporal difference method to update the Q-function and can asymptotically learn the optimal Q-function. Transfer Learning aims to utilize the learned knowledge from source tasks to help new tasks. For supervised learning, it has been shown that transfer learning has the potential to significantly improve the sample complexity of the new tasks. Considering that data collection in RL is both more time and cost consuming and Q-learning converges slowly comparing to supervised learning, different kinds of transfer RL algorithms are designed. However, most of them are heuristic with no theoretical guarantee of the convergence rate. Therefore, it is important for us to clearly understand when and how will transfer learning help RL method and provide the theoretical guarantee for the improvement of the sample complexity. In this paper, we propose to transfer the Q-function learned in the source task to the target in the Q-learning of the new task when certain safe conditions are satisfied. We call this new transfer Q-learning method *target transfer Q-Learning*. The safe conditions are necessary to avoid the harm to the new tasks brought by the transfer target and thus ensure the convergence of the algorithm. We study the convergence rate of the target transfer Q-learning. We prove that if the two tasks are similar with respect to the MDPs, the optimal Q-functions of the two tasks are similar which means the error of the transferred target Q-function in the new task is small. Also, the convergence rate analysis shows that the *target transfer Q-Learning* will converge faster than Q-learning if the error of the transferred target Q-function is smaller than the current Q-function in the new task. Based on our theoretical results and the relationship between the Q error and the Bellman error, we design the safe condition as the Bellman error of the transferred target Q-function is less than the current Q-function. Our experiments are consistent with our theoretical founding and verified the effectiveness of our proposed target transfer Q-learning method.

Introduction

Reinforcement Learning (RL) (Sutton, Barto, and others 1998) technologies are very powerful to learn how to interact with environments and have been successfully applied to variants of important applications, such as robotics, computer games and so on (Kober, Bagnell, and Peters 2013; Mnih et al. 2015; Silver et al. 2016; Bahdanau et al. 2016).

Q-learning (Watkins 1989) is one of the most popular RL algorithms which uses temporal difference method to update the Q-function. To be specific, Q-learning maps the current Q-function to a new Q-function by using Bellman operator and use the difference between these two Q-functions to update the Q-function. Since Bellman operator is a contractive mapping, Q-learning will converge to the optimal Q-function (Jaakkola, Jordan, and Singh 1994). Comparing to supervised learning algorithms, Q-learning converges much slower due to the interactions with the environment. At the same time, the data collection is both very time and cost consuming in RL. Thus, it is crucial for us to utilize available information to save the sample complexity of Q-Learning.

Transfer learning aims to improve the learning performance on a new task by utilizing knowledge/model learned from source tasks. Transfer learning has a long history in supervised learning (Li, Yang, and Xue 2009; Pan, Yang, and others 2010; Oquab et al. 2014). Recently, by leveraging the experiences from supervised transfer learning, researchers developed different kinds of transfer learning methods for RL, which can be categorized into three classes: (1) *instance transfer* in which old data will be reused in the new task (Sunmola and Wyatt 2006; Zhan and Taylor 2015); (2) *representation transfer* such as reward shaping and basis function extraction (Konidaris and Barto 2006; Barreto et al. 2017); (3) *parameter transfer* (Song et al. 2016) in which the parameters of the source task will be partially merged into the model of the new task. While supervised learning is a pure optimization problem, reinforcement learning is a more complex control problem. To the best of our knowledge, most of the existing transfer reinforcement learning algorithms are heuristic with no theoretical guarantee of the convergence rate (Bone 2008), (Taylor and Stone 2009) and (Lazaric 2012). As mentioned by (Spector and Belongie 2017), the transfer learning method potentially do not work or even harm to the new tasks and we do not know the reason since the absence of the theory. Therefore, it is very important

*This work was done when the first author was visiting Microsoft Research Asia.

Copyright © 2019, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

for us to clearly understand how and when transfer learning will help reinforcement learning save sample complexity.

In this paper, we design a novel transfer learning method for Q-learning in RL with theoretical guarantee. Different from the existing transfer RL algorithms, we propose to transfer the Q-function learned in the source task as the temporal difference update target of the new task when certain safe conditions are satisfied. We call this new transfer Q-learning method *target transfer Q-learning*. The intuitive motivation is that when the two RL tasks are similar to each other, their optimal Q-function will be similar which means the transferred target is better (the error is smaller than the current Q-function). Combine it with that a better target Q-function in Q-learning will help to accelerate the convergence, we may expect that the *target transfer Q-learning* method will outperform the Q-learning. The safe conditions are necessary to avoid the harm to the new tasks and thus ensure the convergence of the algorithm.

We prove that target transfer Q-learning has the theoretical guarantee of convergence rate. Furthermore, if the two MDPs and thus the optimal Q-functions in the source and new RL tasks are similar, the target transfer Q-learning converges faster than Q-learning. To be specific, we prove the error of target transfer Q-learning consists of two errors: the initialization error and the sampling error. Both of the errors are increasing with the the product of discount factor γ and the *relative Q-function error ratio* β (*error ratio* for simplicity) which measures the relative error of the target Q-function comparing with the current Q-function in the new task. We called $\gamma\beta$ discounted relative Q-function error ratio(*discounted error ratio* for simplicity). The smaller the discounted error ratio is, the faster the convergence is. And if the discounted error ratio is larger than 1, the convergence will no longer guaranteed.

If the two RL tasks are similar, the learned Q-function in the source task will be close to the optimal Q-function comparing to the current Q-function in the new task. Thus, the discounted error ratio $\gamma\beta$ will be small(especially for the early stage) when we transfer the learned Q-function from the source task to the target of the new task. Please note that the traditional Q-learning is a special case for target transfer Q-learning with constant discounted error ratio γ .

Therefore, our convergence analysis for target transfer Q-learning help us design the safe condition. We can transfer the target if it will lead the discounted error ratio $\gamma\beta$ smaller than 1 . We call it *error ratio safe condition*. Specifically, in the early stage of the training, the Q-function in the new task is not fully trained, the learned Q-function in the source task it a better choice with a smaller error ratio. With the updating of the Q-function in the new task, its error ratio becomes larger. When its discounted error ratio is close or larger than 1, the safe condition will not be satisfied, and we will stop transferring the target to avoid the harm brought by the transfer learning. Following the standard way in Q-learning, we estimate the error ratio about the error of the Q-function w.r.t the optimal Q-function by the Bellman error.

Our experiments on synthetic MDPs fully support our convergence analysis and verify the effectiveness of our proposed target transfer Q-Learning with error ratio safe condition.

Related Work

This section briefly outline related work in transfer learning in reinforcement learning.

Transfer Learning in RL(Taylor and Stone 2009) (Lazarcic 2012) aims to improve learning in new MDP tasks by borrowing knowledge from a related but different learned MDP tasks. In paper (Laroche and Barlier 2017), the authors propose to use instance transfer in the Transfer Reinforcement Learning with Shared Dynamics (TRLSD) setting in which only the reward function is different between MDPs. In paper (Gupta et al. 2017), the authors propose to use the representation transfer and learned the invariant feature space. The papers (Karimpanal and Bouffanais 2018; Song et al. 2016) propose to use the parameter transfer to guide the exploration or to initialize the Q-function of the new task directly. In paper (Al-Shedivat et al. 2017), the authors propose to use the meta-learning method to do transfer learning in RL. All these works are empirically evaluated and no theoretical analysis for the convergence rate.

There are few works that have the convergence analysis. In paper (Barreto et al. 2017), the authors use the representation transfer but only consider the TRLSD setting. (Zhan and Taylor 2015) propose a method by using instance transfer. They gives the theoretical analysis of the asymptotic convergence and no finite sample performance guarantee.

Q Learning Background

Consider the reinforcement learning problem with Markov decision process (MDP) $M \triangleq (\mathcal{S}, \mathcal{A}, P, r, \gamma)$, where \mathcal{S} is the state space, \mathcal{A} is the action space, $P = \{P_{s,s'}^a; s, s' \in \mathcal{S}, a \in \mathcal{A}\}$ is the transition matrix and $P_{s,s'}^a$ is the transition probability from state s to state s' after taking action a , $r = \{r(s, a); s \in \mathcal{S}, a \in \mathcal{A}\}$ is the reward function and $r(s, a)$ is the reward received at state s if taking action a , and $0 < \gamma < 1$ is the discount factor. A policy $\pi : \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ indicates the probability to take each action at each state. Value function for policy π is defined as: $V^\pi(s) \triangleq E[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s, \pi]$. Action value function for policy π is also called Q-function and is defined as:

$$Q^\pi(s, a) \triangleq E \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s, a_0 = a, \pi \right].$$

Without loss of generality, we assume that the rewards all lie between 0 and 1. The optimal policy is denoted π^* and has value function $V_M^*(s)$ and Q value function $Q_M^*(s, a)$.

As we know, the Q-function in RL satisfies the following Bellman equation:

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{\substack{\tilde{a} \sim \pi(a|s) \\ s' \sim P(s'|s, a)}} [Q^\pi(s', \tilde{a}) | s_t = s]$$

Denote the right hand side(RHS) of the equation as $T^\pi Q^\pi(s, a)$, T^π is called Bellman operator for policy π . Similar, consider the optimal Bellman equation:

$$Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{\substack{\tilde{a} \sim \pi^*(a|s) \\ s' \sim P(s'|s, a)}} [Q^*(s', \tilde{a}) | s_t = s]$$

(RHS) of the equation is been denoted as $T^\pi Q^\pi(s, a)$, T^* is called optimal Bellman operator. It can be proved that the optimal Bellman operator is a contraction mapping for the Q-function. We know that there is a unique fix point which is optimal Q-function by contraction mapping theorem. Q-learning algorithm is designed by the above theory. Watkins introduced the Q-learning algorithm to estimate the value of state-action pairs in discounted MDPs (Watkins 1989) :

$$Q_{t+1}(s, a) = (1 - \alpha_t)Q_t(s, a) + \alpha_t \left(r_t(s, a) + \gamma \max_{\tilde{a}} Q_t(s', \tilde{a}) \right)$$

We introduce the max norm error to measure the quality of Q-function:

$$\text{MNE}(Q) = \max_{s, a} |Q(s, a) - Q^*(s, a)|.$$

Target Transfer Q-Learning

First of all, we formalize transfer learning in RL problem. Secondly, We propose our new transfer Q-learning method Target Transfer Q-Learning (TTQL) and introduce the intuition.

Transfer Learning in RL (Taylor and Stone 2009) (Lazarcic 2012) aims to improve learning in new MDP tasks by borrowing knowledge from a related but different learned MDP tasks.

According to the definition of MDPs, $M \triangleq (\mathcal{S}, \mathcal{A}, P, r, \gamma)$, we consider the situation that two MDPs are different in transition probability P , reward function r and discount factor γ . Assume there are two MDPs: source MDP $M_1 = (\mathcal{S}, \mathcal{A}, P_1, r_1, \gamma_1)$ and new MDP $M_2 = (\mathcal{S}, \mathcal{A}, P_2, r_2, \gamma_2)$, Q_1^* and Q_2^* are the corresponding optimal Q-functions. Let M_1 be the source domain and we have already learned the Q_1^* . The goal of transfer in RL considered in this work is how we can use the information of M_1 and Q_1^* to achieve learning speed improvement in M_2 .

To solve the problem mentioned above, we propose to use TTQL method. TTQL use the Q-function learned from the source task as the target Q-function in the new task when safe conditions satisfied. The safe condition ensures that the transferred target only appears if it can help to accelerate the training. Otherwise we will replace it with the current Q-function in the new MDP's learning progress. We describe the TTQL in Algorithm 1.

The intuitive motivation is that when the two RL tasks are similar to each other, their optimal Q-function will be similar. Thus the transferred target is better (the error is smaller than the current Q-function) and the better target can help to accelerate the convergence.

We define the distance between two MDPs as $\Delta(M_1, M_2)$

$$\Delta(M_1, M_2) = \max_{s, a} |Q_1^*(s, a) - Q_2^*(s, a)|.$$

The following Proposition 1 shows the relation between the distance of two MDPs and the component of two MDPs.

Proposition 1. Assume two MDPs, $M_1 = (\mathcal{S}, \mathcal{A}, P_1, r_1, \gamma_1)$ and $M_2 = (\mathcal{S}, \mathcal{A}, P_2, r_2, \gamma_2)$, Let the corresponding optimal

Algorithm 1 Target Transfer Q Learning

Require: initial Q-learning Q_1 , source task learned Q-learning Q_{source}^* , total step n

- 1: **for** $t = 1, \dots, n$ **do**
- 2: $\alpha_t = \frac{1}{t+1}$
- 3: $\text{flag} = \text{safe-condition}(Q_{source}^*, Q_t(\cdot, \cdot))$
- 4: **if** $\text{flag} = \text{True}$ **then**
- 5: $Q_{target} = Q_{source}^*$
- 6: **else**
- 7: $Q_{target} = Q_t$
- 8: **end if**
- 9: **for** $s \in \mathcal{S}, a \in \mathcal{A}$ **do**
- 10: $Q_{t+1}(s, a) = (1 - \frac{1}{n})Q_t(s, a) + \frac{1}{n}(r(s, a) + \gamma \max_{\tilde{a}} Q_{target}(s', \tilde{a}))$
- 11: **end for**
- 12: **end for**

Ensure: Q_{n+1}

Q-functions be Q_1^* and Q_2^* , then we have

$$\begin{aligned} \Delta(M_1, M_2) &= \|Q_1^* - Q_2^*\|_\infty \leq \tilde{\Delta}(M_1, M_2) \quad (1) \\ &\triangleq \frac{\|r_1 - r_2\|_\infty}{1 - \gamma'} + \frac{\gamma'' \|r'\|_\infty}{(1 - \gamma'')^2} \|P_1 - P_2\|_\infty + \frac{|\gamma_1 - \gamma_2|}{(1 - \gamma_1)(1 - \gamma_2)} \|r''\|_\infty. \end{aligned}$$

for $\forall(\gamma', \gamma'', r', r'') \in \Omega$, where Ω is the available combination of the $(\gamma_1, \gamma_2, \gamma_1, \gamma_2)$.

Proof. Without loss of generality, we assume $\gamma_1 \leq \gamma_2$, $\|r_2\|_\infty \leq \|r_1\|_\infty$, we will show that other cases can be proved similarly. We define the following auxiliary MDPs: $\hat{M}_3 = (\mathcal{S}, \mathcal{A}, P_1, r_2, \gamma_1)$, $\hat{M}_4 = (\mathcal{S}, \mathcal{A}, P_2, r_2, \gamma_1)$, and let the corresponding optimal Q-functions be Q_3^* and Q_4^* . We have

$$\|Q_1^* - Q_2^*\|_\infty \quad (2)$$

$$= \|Q_1^* - Q_3^* + Q_3^* - Q_4^* + Q_4^* - Q_2^*\|_\infty \quad (3)$$

$$\leq \|Q_1^* - Q_3^*\|_\infty + \|Q_3^* - Q_4^*\|_\infty + \|Q_4^* - Q_2^*\|_\infty \quad (4)$$

Notice that in each term, two MDPs are only different in one component. Using the results of (Csaji and Monostori 2008), we have that $\|Q_1^* - Q_3^*\|_\infty \leq \frac{\|r_1 - r_2\|_\infty}{1 - \gamma_1}$, $\|Q_3^* - Q_4^*\|_\infty \leq \frac{\gamma_1 \|r_2\|_\infty}{(1 - \gamma_1)^2} \|P_1 - P_2\|_\infty$, $\|Q_4^* - Q_2^*\|_\infty \leq \frac{|\gamma_1 - \gamma_2|}{(1 - \gamma_1)(1 - \gamma_2)} \|r_2\|_\infty$. Combine the above upper bound and set $\gamma' = \gamma_1, \gamma'' = \gamma_1, r' = r_2, r'' = r_2$, we can get the in-equation (1).

In other situation, we can construct auxiliary MDPs like above and use the similar procedure to prove the theorem. After traversing all the available combination of the $(\gamma_1, \gamma_2, \gamma_1, \gamma_2)$, we can prove the Proposition 1 □

By the Proposition 1, we can conclude that if the two RL tasks are similar, in the sense of that the component of two MDPs are similar, the learned Q-function in the source task will be close to the optimal Q in the new task.

A question is that when to transfer the target will have performance guarantee. Here, we need safe conditions which

are necessary to avoid the harm to the new tasks and thus ensure the convergence of the algorithm. We can now heuristically relate it to the distance between two MDPs and the current learning quality. The concrete value of the safe condition need to further investigate through quantified theoretical analysis and we present these result in the following section.

Convergence Rate of TTQL

In this section, we present the convergence rate of the Target Transfer Q Learning (TTQL) and make discussions for the key factor that influence the convergence. Theorem 1 analysis the convergence of the target transfer Q learning. Theorem 2 and 3 analysis two key factors of the convergence rate. Theorem 4 discuss the convergence rate for the TTQL totally.

First of all, Theorem 1 analysis the convergence rate for the target transfer method which is

$$Q_{t+1}(s, a) = (1 - \frac{1}{n})Q_t(s, a) + \frac{1}{n} \left(r(s, a) + \gamma \max_{\tilde{a}} Q_{target}(s', \tilde{a}) \right)$$

For simplicity, we denote $E_n = \text{MNE}(Q_n)$. We denote the error ratio $\beta_n = \frac{\text{MNE}(Q_{target})}{E_n}$ and β if we do not specify the learning steps n .

Theorem 1. we denote $w_k(\beta_{n-k:n}) = \frac{\prod_{i=n-k}^{n-1} (i + \gamma\beta_i)}{\prod_{i=n-k}^n i}$, $\alpha_n = \frac{\prod_{i=1}^{n-1} (i + \gamma\beta_i)}{\prod_{i=2}^n i}$. If $0 \leq \beta_n \leq 1$, then with probability $1 - \delta$ we have

$$E_n \leq \underbrace{\alpha_n E_1}_{\text{initialization error}} + \underbrace{\sqrt{\frac{\ln 1/\delta \sum_{k=0}^{n-1} w_k^2(\beta_{n-k:n})}{2}}}_{\text{sampling error}}$$

Before showing the proof of Theorem 1, we first introduce a modified Hoeffding inequality lemma which bounds the distance between the weighted sum of the bounded random variable and its expectation.

Lemma 1. Let $a < x_i < b$ almost surely, $S_n = \sum_{i=1}^n w_i x_i$, then we have

$$S_n - E[S_n] \leq \sqrt{\frac{1}{2} \log \frac{1}{\delta} \sum_{k=1}^n w_k^2 (b-a)^2}. \quad (5)$$

Proof. We first prove the inequality $\mathbb{P}(S_n - E[S_n] \geq \epsilon) \leq \exp\left(-\frac{2\epsilon^2}{\sum_{k=1}^n w_k^2 (b-a)^2}\right)$

For $s, \epsilon \geq 0$, Markov's inequality and the independence of

x_i implies

$$\mathbb{P}(S_n - E[S_n] \geq \epsilon) \quad (6)$$

$$= \mathbb{P}\left(e^{s(S_n - E[S_n])} \geq e^{s\epsilon}\right) \quad (7)$$

$$\leq e^{-s\epsilon} \mathbb{E}\left[e^{s(S_n - E[S_n])}\right] \quad (8)$$

$$= e^{-s\epsilon} \mathbb{E}\left[e^{s(\sum_{i=1}^n w_i x_i - E[\sum_{i=1}^n w_i x_i])}\right] \quad (9)$$

$$= e^{-s\epsilon} \prod_{i=1}^n \mathbb{E}\left[e^{s w_i (x_i - E[x_i])}\right] \quad (10)$$

$$\leq e^{-s\epsilon} \prod_{i=1}^n e^{\frac{s^2 w_i^2 (b-a)^2}{8}} \quad (11)$$

$$= \exp\left(-s\epsilon + \frac{1}{8} s^2 (b-a)^2 \sum_{i=1}^n w_i^2\right). \quad (12)$$

Now we consider the minimum of the right hand side of the last inequality as a function of s , and denote

$$g(s) = -s\epsilon + \frac{1}{8} s^2 (b-a)^2 \sum_{i=1}^n w_i^2$$

Note that g is a quadratic function and achieves its minimum at $s = \frac{4\epsilon}{(b-a)^2 \sum_{i=1}^n w_i^2}$, Thus we get

$$\mathbb{P}(S_n - E[S_n] \geq \epsilon) \leq \exp\left(-\frac{2\epsilon^2}{\sum_{k=1}^n w_k^2 (b-a)^2}\right) \quad (13)$$

We can easily obtain the second part of the Lemma 1 by inverse the inequality. \square

Proof of Theorem 1. Our analysis are derived based on the following synchronous generalized Q-learning setting. Compare with the traditional synchronous Q-learning¹, we replace the target Q-function as the independent Q-function $Q'(s, a)$ rather than the current one $Q_n(s, a)$.

$$\forall s, a : Q_0(s, a) = q(s, a)$$

$$\forall s, a : Q_n(s, a) =$$

$$\left(\frac{n-1}{n}\right) Q_{n-1}(s, a) + \frac{1}{n} \left(r(s, a) + \gamma \max_{\tilde{a}} Q'_{n-1}(s', \tilde{a}) \right) \quad (14)$$

Let $Q'_n(s, a)$ satisfied the following condition ,

$$0 \leq \frac{\max_{s,a} (Q'_n(s, a) - Q^*(s, a))}{\max_{s,a} (Q_n(s, a) - Q^*(s, a))} \leq 1 \quad (15)$$

Note that if we set $Q'_n(s, a) = Q^*_{source}$, we can verify $0 \leq \beta_n \leq 1$ according to inequality 15. First of all, we decompose the update role,

$$\begin{aligned} Q_n(s, a) &= \frac{n-1}{n} Q_{n-1}(s, a) + \frac{1}{n} \left[r(s, a) + \gamma \max_{\tilde{a}} Q'_{n-1}(s', \tilde{a}) \right] \\ &= \frac{n-1}{n} Q_{n-1}(s, a) + \frac{1}{n} \left[r(s, a) + \gamma \max_{\tilde{a}} Q^*(s', \tilde{a}) \right. \\ &\quad \left. + \gamma \max_{\tilde{a}} Q'_{n-1}(s', \tilde{a}) - \gamma \max_{\tilde{a}} Q^*(s', \tilde{a}) \right] \end{aligned}$$

¹It is the same as the commonly used setting or more general(Asadi and Littman 2017), (Even-Dar and Mansour 2003), (Azar et al. 2013) (Haarnoja et al. 2017)).

If we denote $\epsilon_n(s, a) = Q_n(s, a) - Q^*(s, a)$, $x(s') = \gamma \max_{\tilde{a}} Q^*(s', \tilde{a})$ and recall the definition of β_n we can have

$$\begin{aligned} & \epsilon_n(s, a) \\ & \leq \frac{n-1}{n} \epsilon_{n-1}(s, a) + \frac{1}{n} [x(s') - \mathbb{E}_{s'} x(s')] + \frac{1}{n} \gamma \beta_n \epsilon_{n-1}(s', \tilde{a}) \\ & \leq \frac{n-1}{n} \epsilon_{n-1}(s, a) + \frac{1}{n} [x(s') - \mathbb{E}_{s'} x(s')] + \frac{1}{n} \gamma \beta_n E_{n-1} \end{aligned}$$

The last step is right because $\epsilon_n(s, a) \leq E_n$ for $\forall s, a$. Taking maximization of the both sides(RHS) of the inequality and using recursion of E we can have

$$\begin{aligned} E_n & \leq \frac{n-1 + \gamma \beta_n}{n} E_{n-1} + \frac{1}{n} [x(s') - \mathbb{E}_{s'} x(s')] \\ & \leq \frac{\prod_{i=1}^{n-1} (i + \gamma \beta_i)}{\prod_{i=2}^n i} E_1 + \sum_{k=1}^{n-1} \frac{\prod_{i=n-k}^{n-1} (i + \gamma \beta_i)}{\prod_{i=n-k}^n i} [x(s'_k) - \mathbb{E}_{s'} x(s')] \\ & = \alpha_n E_1 + \sum_{k=1}^{n-1} w_k(\beta) [x(s'_k) - \mathbb{E}_{s'} x(s')] \end{aligned}$$

According to Lemma 1(weighted Hoeffding inequality), with probability $1-\delta$, we have

$$E_n \leq \alpha_n E_1 + \sqrt{\frac{\ln 1/\delta \sum_{k=0}^{n-1} w_k^2(\beta_{n-k:n})}{2}} \quad (16)$$

□

The convergence result reveals the how the error ratio β influence the convergence rate. In short, if we can find a better target Q-function, we can learn much more faster.

We can see from the Theorem 1 that there are two key factors that influence the convergence rate. One is the initialization error $\alpha_n E_1$, the other one is the sampling error $\sqrt{\frac{\ln 1/\delta \sum_{k=0}^{n-1} w_k^2(\beta_{n-k:n})}{2}}$. To make it clear, we analysis the order of these two terms in 2 and 3 respectively.

Theorem 2. Denote $w_k(\beta_{n-k:n}) = \frac{\prod_{i=n-k}^{n-1} (i + \gamma \beta_i)}{\prod_{i=n-k}^n i}$, and $\beta_i \leq \beta^*$ for $\forall i \leq n$, we have

$$\sum_{k=0}^{n-1} (w_k(\beta_{n-k:n}))^2 \leq \begin{cases} \frac{e^{2\gamma\beta^*}}{n^{2-2\gamma\beta^*}} \left(\frac{n^{1-2\gamma\beta^*}}{1-2\gamma\beta^*} - \frac{1}{1-2\gamma\beta^*} + 1 \right), & \gamma\beta^* \neq 0.5 \\ \frac{(n-2)^{2\gamma\beta^*}}{n^2} e^{2\gamma\beta^*} (1 + \ln(n)), & \gamma\beta^* = 0.5 \end{cases}$$

Based on the results of Theorem 2, we can get the following corollary directly.

Corollary 1. The order of $\sum_{k=0}^{n-1} (w_k(\beta_{n-k:n}))^2$ is:

$\mathcal{O}(\frac{1}{n})$, if $\gamma\beta^* < 0.5$.

$\mathcal{O}(\frac{1}{n^{2-2\gamma\beta^*}})$, if $0.5 < \gamma\beta^* < 1$.

$\mathcal{O}(\frac{1}{n^{2-2\gamma\beta^*}} \ln(n))$, if $\gamma\beta^* = 0.5$.

The sufficient condition for the $\lim_{n \rightarrow \infty} \sum_{k=0}^{n-1} (w_k(\beta^*))^2 = 0$ is $\gamma\beta^* < 1$

Before showing the proof of Theorem 2, we first introduce a Lemma which will be used.

Lemma 2. If $a < b$, $\sum_{i=a}^b \frac{1}{i} \leq \frac{1}{a} + \ln(b) - \ln(a)$.

Proof.

$$\begin{aligned} \sum_{i=a}^b \frac{1}{i} & \leq \frac{1}{a} + \sum_{i=a+1}^b \frac{1}{i} \leq \frac{1}{a} + \sum_{i=a+1}^b \int_{k=i-1}^i \frac{1}{k} dk \\ & \leq \frac{1}{a} + \int_{k=a}^b \frac{1}{k} dk \leq \frac{1}{a} + \ln(b) - \ln(a) \end{aligned}$$

□

Proof of Theorem 2.

$$\begin{aligned} \sum_{k=0}^{n-1} (w_k(\beta_{n-k:n}))^2 & \leq \sum_{k=0}^{n-1} \left(\frac{\prod_{i=n-k}^{n-1} (i + \gamma\beta^*)}{\prod_{i=n-k}^n i} \right)^2 \\ & \stackrel{(a)}{\leq} \sum_{k=0}^{n-1} \exp \left\{ 2 \left[\sum_{i=n-k}^{n-1} \ln(i + \gamma\beta^*) - \sum_{i=n-k}^n \ln i \right] \right\} \\ & \stackrel{(b)}{\leq} \frac{1}{n^2} \sum_{k=0}^{n-1} \exp \left\{ 2 \sum_{i=n-k}^{n-1} [\ln(i + \gamma\beta^*) - \ln i] \right\} \\ & \stackrel{(c)}{\leq} \frac{1}{n^2} \sum_{k=0}^{n-1} \exp \left\{ 2 \sum_{i=n-k}^{n-1} \frac{\gamma\beta^*}{i} \right\} \\ & \stackrel{(d)}{\leq} \frac{1}{n^2} \sum_{k=0}^{n-1} \exp \{ 2\gamma\beta^* [\ln(n-2) - \ln(n-k) + 1] \} \\ & = \frac{(n-2)^{2\gamma\beta^*}}{n^2} e^{2\gamma\beta^*} \sum_{k=0}^{n-1} \frac{1}{(n-k)^{2\gamma\beta^*}} \\ & = \frac{(n-2)^{2\gamma\beta^*}}{n^2} e^{2\gamma\beta^*} \sum_{t=1}^n \frac{1}{t^{2\gamma\beta^*}} \quad (17) \end{aligned}$$

We rewrite the product term in (a) into the summarization term. Then we drop one term outside of the summarization to align the i sum from $n-k$ to $n-1$ in (b). (c) follows the concave property of the \ln function. (d) follows the relation between summarization and integral as shown in Lemma 2. The last two terms is right because we only rearrange the term and write it simply.

If $\gamma\beta^* = 0.5$, $2\gamma\beta^* = 1$,

$$\sum_{k=0}^{n-1} (w_k(\beta_{n-k:n}))^2 \leq \frac{1}{n^{2-2\gamma\beta^*}} e^{\frac{2\gamma\beta^*}{n-1}} (1 + \ln(n))$$

If $\gamma\beta^* \neq 0.5$,

$$\sum_{k=0}^{n-1} (w_k(\beta^*))^2 \leq \underbrace{\frac{1}{n^{2-2\gamma\beta^*}}}_{(e)} \underbrace{e^{2\gamma\beta^*}}_{(f)} \left(\underbrace{\frac{n^{1-2\gamma\beta^*}}{1-2\gamma\beta^*}}_{(g)} - \underbrace{\frac{1}{1-2\gamma\beta^*} + 1}_{(h)} \right)$$

Note that term (f) is a constant.

If $\gamma\beta^* < 0.5$, term(g) will dominant the order, $\sum_{k=0}^{n-1} (w_k(\beta_{n-k:n}))^2$ will be $\mathcal{O}(\frac{1}{n})$.
If $\gamma\beta^* > 0.5$, term(h) will dominant the order,

$\sum_{k=0}^{n-1} (w_k(\beta_{n-k:n}))^2$ will be $\mathcal{O}(\frac{1}{n^{2-2\gamma\beta^*}})$.
If $\gamma\beta^* = 0.5$, $\sum_{k=0}^{n-1} (w_k(\beta_{n-k:n}))^2$ will be $\mathcal{O}(\frac{1}{n^{2-2\gamma\beta^*}} \ln(n))$.
In all case, the (18) will converge to 0 as n will go to ∞ . \square

Note that if $\gamma\beta^* < 1$. The theorem 2 shows, $\sum_{k=0}^{n-1} w_k^2$ converges to 0 and the convergence rate is highly related to the $\gamma\beta^*$. The next theorem shows the upper bound of the coefficient α_n in initialization error.

Theorem 3. Denote $\alpha_n = \frac{\prod_{i=1}^{n-1} (i + \gamma\beta^*)}{\prod_{i=2}^n i}$, and $\beta_i \leq \beta^*$ for $\forall i \leq n$, we can bound α_n as:

$$\alpha_n \leq \frac{(n-1)^{\gamma\beta^*}}{n} (1 + \gamma\beta^*) e^{(0.5 - \ln 2)\gamma\beta^*} = \frac{C_{\gamma, \beta^*}^1}{n^{1-\gamma\beta^*}}. \quad (18)$$

where $C_{\gamma, \beta^*}^1 = (1 + \gamma\beta^*) e^{(0.5 - \ln 2)\gamma\beta^*}$ is a constant

Proof of Theorem 3.

$$\alpha_n \leq \frac{\prod_{i=1}^{n-1} (i + \gamma\beta^*)}{\prod_{i=2}^n i} \quad (19)$$

$$= \exp \left\{ \sum_{i=1}^{n-1} \ln(i + \gamma\beta^*) - \sum_{i=2}^n \ln i \right\} \quad (20)$$

$$= (1 + \gamma\beta^*) \exp \left\{ \sum_{i=2}^{n-1} (\ln(i + \gamma\beta^*) - \ln i) - \ln n \right\} \quad (21)$$

$$\leq (1 + \gamma\beta^*) \exp \left\{ \sum_{i=2}^{n-1} \left(\frac{\gamma\beta^*}{i} \right) - \ln n \right\} \quad (22)$$

$$\leq (1 + \gamma\beta^*) \exp \{ \gamma\beta^* (0.5 + \ln(n-1) - \ln 2) - \ln n \} \quad (23)$$

$$\leq \frac{(n-1)^{\gamma\beta^*}}{n} (1 + \gamma\beta^*) e^{(0.5 - \ln 2)\gamma\beta^*} \quad (24)$$

We rewrite the product term in the second equation into the summarization term. The third equation is rearrange the terms. The first inequality follows the concave property of \ln function. The second inequality follows the relation between summarization and integral(Lemma 2). \square

Note that if $\gamma\beta^* < 1$. The theorem 3 shows, α_n converge to 0 and the convergence rate is in order $\mathcal{O}(\frac{1}{n^{1-\gamma\beta^*}})$.

Combining Theorem 1, 2 and 3, we have the following Theorem:

Theorem 4. The TTQL will converge if we set the *safe condition* as

$$\hat{\beta}_n = \frac{\Delta(M_1, M_2)}{E_n} \leq 1.$$

And the convergence rate is:

$$E_n \leq \begin{cases} \mathcal{O}(\frac{1}{n^{1-\gamma\beta}} E_1 + \sqrt{\frac{1}{n}}), & \text{if } \gamma\beta < 0.5 \\ \mathcal{O}(\frac{1}{n^{1-\gamma\beta}} E_1 + \frac{1}{n^{1-\gamma\beta}} \sqrt{\ln n}), & \text{if } \gamma\beta = 0.5 \\ \mathcal{O}(\frac{1}{n^{1-\gamma\beta}} E_1 + \frac{1}{n^{1-\gamma\beta}}), & \text{if } 0.5 < \gamma\beta < 1 \end{cases}. \quad (25)$$

Note that if the safe condition is satisfied, we set $Q_{target} = Q_{source}^*$ and β

We would like to make the following discussion:

(1) The distance between two MDPs influence the convergence rate. According to the Proposition 1, if two MPDs have the similar components(P, r, γ), the optimal Q-function of these two MDPs will be closed. The discounted error ratio $\gamma\beta_n$ will be relatively small in this situation and the convergence rate will be improved.

(2) Q-learning is the special case. Please note that the traditional Q-learning is a special case for target transfer Q-learning with $Q_{target} = Q_{n-1}$. Thus the error ratio is a constant and $\beta_n = 1$ and our results reduce to the previous (Szepesvári 1998). It shows that if the $\beta < 1$ in TTQL, the TTQL converge faster than traditional Q-learning.

(3) The TTQL method do converge with the safe condition. As shown in Theorem 4, the TTQL method will converge. And the convergence rate changes under different discounted error ratio $\gamma\beta$. The smaller $\gamma\beta$ will lead to a quicker convergence rate. Intuitively, smaller β means that Q' provides more information about the optimal Q-function. Besides, the discount factor γ can be viewed as the "horizon" of the infinite MDPs. Smaller γ means that the expected long-term return is less influenced by the future information and the immediate reward is assigned more weights.

(4) Safe condition is necessary. As mentioned above, the safe condition is defined as $\hat{\beta}_n \leq 1$. If the safe condition is satisfied, we set $Q_{target} = Q_{source}^*$ and $\gamma\beta_n = \gamma\hat{\beta}_n \leq \gamma < 1$. If safe condition is not satisfied, we set $Q_{target} = Q_n$ and $\gamma\beta_n = \gamma < 1$. So with the safe condition, TTQL algorithms do converge at any situation. At the beginning of the new task training, due to the large error of the current Q-function, $\beta_n = \hat{\beta}_n$ will be relatively small and the transfer learning will be greatly helpful. Speedup would come down as the error of current Q-function, become smaller. Finally when β is equal to or larger than one we need to remove the transfer Q target which means to set $\beta = 1$ to avoid the harm brought by the transfer learning.

Discussion for Error Ratio Safe Condition

Until now, we can conclude that TTQL will converge. TTQL method need the safe condition to guarantee the convergence. In this section, We discuss the safe conditions.

At the beginning, we propose the safe condition is that can guarantee the algorithms convergence generally. Heuristically, the safe condition is related to the distance between two MDPs and the quality of the current value function. Then according to the Theorem 1, we know that the safe condition is $\hat{\beta}_n \leq 1$ which we called error ratio safe condition. Under the transfer learning in RL setting, it means that the distance between two MDPs need to be smaller than the error of the current Q-function. In the real algorithms, it is impossible to calculate the error of the current Q-function $MNE(Q_n)$ and the distance between two MDPs precisely. However it is easy to calculate the bellman error $MINBE(Q(s, a)) = \max_{s,a} |Q(s, a) - (r(s, a) + \gamma E_{s'} \max_{\tilde{a}} (Q(s', \tilde{a})))|$. We

can prove that these two metrics follow the relationship as:

$$\text{MNE}(Q) \leq \frac{\text{MNBE}(Q)}{1 - \gamma}.$$

Following the standard way in Q-learning, we estimate the error ratio about the error of the Q-function w.r.t the optimal Q-function by the Bellman error.

Algorithm 2 Error Ratio Safe Condition

Require: learned Q_1^* , current Q-function Q_n

- 1: **if** $\text{MNBE}(Q_1^*) \leq \text{MNBE}(Q_n)$ **then**
- 2: flag = True
- 3: **else**
- 4: flag = False
- 5: **end if**

Ensure: flag

Proof of the relation between MNE and MNBE. Denote $\mathcal{B}Q(s, a) = r(s, a) - \gamma \mathbb{E}_{s'} \max_{\tilde{a}} Q(s', \tilde{a})$ as bellman operator.

$$\begin{aligned} & \text{MNE}(Q) \\ & \leq \|Q(s, a) - \mathcal{B}Q^*(s, a)\|_\infty + \|\mathcal{B}Q^*(s, a) - Q^*(s, a)\|_\infty \\ & \leq \text{MNBE}(Q) + \|\gamma \mathbb{E}_{s'} \max_{\tilde{a}} Q(s', \tilde{a}) - \gamma \mathbb{E}_{s'} \max_{\tilde{a}} Q^*(s', \tilde{a})\| \\ & \leq \text{MNBE}(Q) + \gamma \text{MNE}(Q) \end{aligned}$$

So we can proof that

$$\text{MNE}(Q) \leq \frac{\text{MNBE}(Q)}{1 - \gamma}.$$

□

Experiment

In this section, we report our simulation experiments to support our convergence analysis and verified the effectiveness of our proposed target transfer Q-Learning with the error ratio safe condition.

We consider the general MDP setting. We construct the random MDP by generating the transition probability $P(s'|s, a)$, reward function $r(s, a)$ and discount factor γ and fixing the state and action space size as 50.

First of all, we generate 9 different MDPs ($M_{11} \sim M_{33}$) as source tasks and then generate the new MDP M_0 . Let M_{11}, M_{12}, M_{13} be different from M_0 in γ and the distance from M_1 . and M_0 increase as $M_{11} < M_{12} < M_{13}$. Similarly, MDPs M_{21}, M_{22}, M_{23} is different from M_0 in r , and MDPs M_{31}, M_{32}, M_{33} is different from M_0 in P . Then we run our algorithm to transfer the Q-function learned on these 9 source MPDs to the new MDP M_0 . The result is shown in Figure 1a, 1b and 1c. Note that the dash line Q is the Q-learning algorithm with no transfer learning, and the solid line with various markers are the TTQL algorithm.

Secondly, we design three MDPs M_4, M_5, M_6 as source task MDPs, and the distance between these MDPs and the target becomes larger and larger. Then we use TTQL to transfer the Q-function learning from them to new MDP M_0 with and without the safe condition. The results is shown in Figure 1d, 1e and 1f. Note that $W - SC$ means that the experiment is

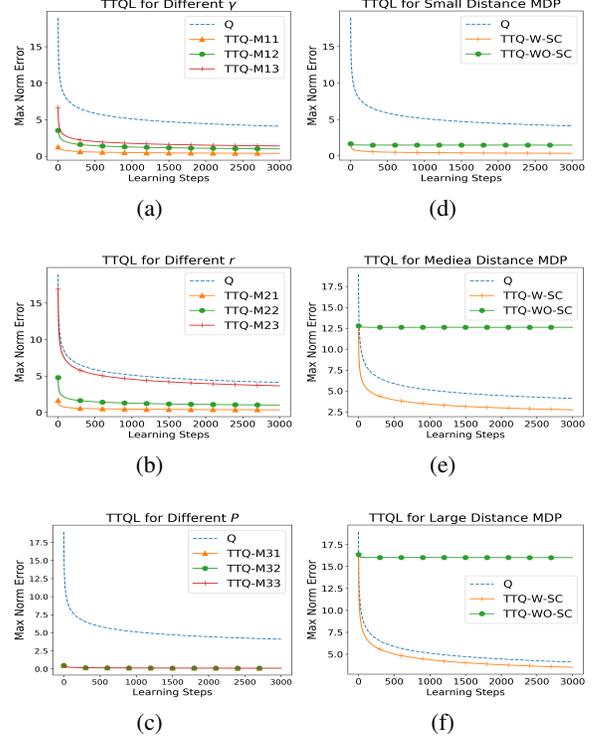


Figure 1: Left three figures are the learning errors w.r.t the three types of different MDPs (Be different in γ , r , P respectively). Right three figures are the learning error w.r.t the three different distance transfer task and both training with/without the safe condition.

run with the safe condition and $WO - SC$ means without the safe condition.

We have the following observations. (1) TTQL method outperforms Q-learning in all experiments. (2) Running TTQL on the more similar MDPs will lead to the faster convergence rate. Note that the curve in Figure 1e are closed to each other. It is because the infinity norm of the P will be small because the scale of the P is small and is consistent with the Proposition 1. (3) The safe condition is necessary to ensure the convergence of the algorithms in various situation. All these observations are consistent with our theoretical findings.

Conclusion

In this paper, we proposed a new transfer learning in RL method *target transfer Q-learning*(TTQL). The method transfer the Q-function learned in the source task to the target of Q-learning in the new task when the safe conditions are satisfied. We prove the TTQL method do converge with the safe condition and the convergence rate is quicker than Q-learning if the two MDPs are not faraway from each other. The theoretical analysis helps to design safe conditions which is key to guarantee the convergence of TTQL. As far as we known, it is the first convergence rate guaranteed transfer leaning in reinforcement learning algorithm. In the future, we will apply

the TTQL to the more complex tasks and study convergence rate for the TTQL with complex function approximation such as the neural network.

References

- [Al-Shedivat et al. 2017] Al-Shedivat, M.; Bansal, T.; Burda, Y.; Sutskever, I.; Mordatch, I.; and Abbeel, P. 2017. Continuous adaptation via meta-learning in nonstationary and competitive environments. *arXiv preprint arXiv:1710.03641*.
- [Asadi and Littman 2017] Asadi, K., and Littman, M. L. 2017. An alternative softmax operator for reinforcement learning. In Precup, D., and Teh, Y. W., eds., *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, 243–252. International Convention Centre, Sydney, Australia: PMLR.
- [Azar et al. 2013] Azar, M.; Munos, R.; Ghavamzadeh, M.; and Kappen, H. 2013. Speedy q-learning: a computationally efficient reinforcement learning algorithm with a near optimal rate of convergence. *Journal of Machine Learning Research* 1–26.
- [Bahdanau et al. 2016] Bahdanau, D.; Brakel, P.; Xu, K.; Goyal, A.; Lowe, R.; Pineau, J.; Courville, A.; and Bengio, Y. 2016. An actor-critic algorithm for sequence prediction. *arXiv preprint arXiv:1607.07086*.
- [Barreto et al. 2017] Barreto, A.; Munos, R.; Schaul, T.; and Silver, D. 2017. Successor features for transfer in reinforcement learning. In *Advances in neural information processing systems*.
- [Bone 2008] Bone, N. 2008. A survey of transfer learning methods for reinforcement learning.
- [Csáji and Monostori 2008] Csáji, B. C., and Monostori, L. 2008. Value function based reinforcement learning in changing markovian environments. *Journal of Machine Learning Research* 9(Aug):1679–1709.
- [Even-Dar and Mansour 2003] Even-Dar, E., and Mansour, Y. 2003. Learning rates for q-learning. *Journal of Machine Learning Research* 5(Dec):1–25.
- [Gupta et al. 2017] Gupta, A.; Devin, C.; Liu, Y.; Abbeel, P.; and Levine, S. 2017. Learning invariant feature spaces to transfer skills with reinforcement learning. *arXiv preprint arXiv:1703.02949*.
- [Haarnoja et al. 2017] Haarnoja, T.; Tang, H.; Abbeel, P.; and Levine, S. 2017. Reinforcement learning with deep energy-based policies. In *International Conference on Machine Learning*, 1352–1361.
- [Jaakkola, Jordan, and Singh 1994] Jaakkola, T.; Jordan, M. I.; and Singh, S. P. 1994. Convergence of stochastic iterative dynamic programming algorithms. In *Advances in neural information processing systems*, 703–710.
- [Karimpanal and Bouffanais 2018] Karimpanal, T. G., and Bouffanais, R. 2018. Self-organizing maps as a storage and transfer mechanism in reinforcement learning. *CoRR* abs/1807.07530.
- [Kober, Bagnell, and Peters 2013] Kober, J.; Bagnell, J. A.; and Peters, J. 2013. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research* 32(11):1238–1274.
- [Konidaris and Barto 2006] Konidaris, G., and Barto, A. 2006. Autonomous shaping: Knowledge transfer in reinforcement learning. In *Proceedings of the 23rd international conference on Machine learning*, 489–496. ACM.
- [Laroche and Barlier 2017] Laroche, R., and Barlier, M. 2017. Transfer reinforcement learning with shared dynamics. In *AAAI*.
- [Lazaric 2012] Lazaric, A. 2012. Transfer in reinforcement learning: a framework and a survey. In *Reinforcement Learning*. Springer. 143–173.
- [Li, Yang, and Xue 2009] Li, B.; Yang, Q.; and Xue, X. 2009. Transfer learning for collaborative filtering via a rating-matrix generative model. In *Proceedings of the 26th annual international conference on machine learning*, 617–624. ACM.
- [Mnih et al. 2015] Mnih, V.; Kavukcuoglu, K.; Silver, D.; Rusu, A. A.; Veness, J.; Bellemare, M. G.; Graves, A.; Riedmiller, M.; Fidjeland, A. K.; Ostrovski, G.; et al. 2015. Human-level control through deep reinforcement learning. *Nature* 518(7540):529–533.
- [Oquab et al. 2014] Oquab, M.; Bottou, L.; Laptev, I.; and Sivic, J. 2014. Learning and transferring mid-level image representations using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 1717–1724.
- [Pan, Yang, and others 2010] Pan, S. J.; Yang, Q.; et al. 2010. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering* 22(10):1345–1359.
- [Silver et al. 2016] Silver, D.; Huang, A.; Maddison, C. J.; Guez, A.; Sifre, L.; Van Den Driessche, G.; Schrittwieser, J.; Antonoglou, I.; Panneershelvam, V.; Lanctot, M.; et al. 2016. Mastering the game of go with deep neural networks and tree search. *Nature* 529(7587):484–489.
- [Song et al. 2016] Song, J.; Gao, Y.; Wang, H.; and An, B. 2016. Measuring the distance between finite markov decision processes. In *Proceedings of the 2016 International Conference on Autonomous Agents & Multiagent Systems*, 468–476. International Foundation for Autonomous Agents and Multiagent Systems.
- [Spector and Belongie 2017] Spector, B., and Belongie, S. 2017. Sample-efficient reinforcement learning through transfer and architectural priors. In *Advances in neural information processing systems*.
- [Sunmola and Wyatt 2006] Sunmola, F. T., and Wyatt, J. L. 2006. Model transfer for markov decision tasks via parameter matching. In *Proceedings of the 25th Workshop of the UK Planning and Scheduling Special Interest Group (PlanSIG 2006)*.
- [Sutton, Barto, and others 1998] Sutton, R. S.; Barto, A. G.; et al. 1998. *Reinforcement learning: An introduction*. MIT press.
- [Szepesvári 1998] Szepesvári, C. 1998. The asymptotic convergence-rate of q-learning. In *Advances in Neural Information Processing Systems*, 1064–1070.

- [Taylor and Stone 2009] Taylor, M. E., and Stone, P. 2009. Transfer learning for reinforcement learning domains: A survey. *Journal of Machine Learning Research* 10:1633–1685.
- [Watkins 1989] Watkins, C. J. C. H. 1989. *Learning from delayed rewards*. Ph.D. Dissertation, King’s College, Cambridge.
- [Zhan and Taylor 2015] Zhan, Y., and Taylor, M. E. 2015. Online transfer learning in reinforcement learning domains. *arXiv preprint arXiv:1507.00436*.