

On Convergence of Average-Reward Off-Policy Control Algorithms in Weakly Communicating MDPs

Yi Wan

University of Alberta

Richard S. Sutton

University of Alberta, DeepMind

WAN6@UALBERTA.CA

RSUTTON@UALBERTA.CA

Abstract

We show two average-reward off-policy control algorithms, Differential Q-learning (Wan, Naik, & Sutton 2021a) and RVI Q-learning (Abounadi Bertsekas & Borkar 2001), converge in weakly communicating MDPs. Weakly communicating MDPs are the most general MDPs that can be solved by a learning algorithm with a single stream of experience. The original convergence proofs of the two algorithms require that the solution set of the average-reward optimality equation only has one degree of freedom, which is not necessarily true for weakly communicating MDPs. To the best of our knowledge, our results are the first showing average-reward off-policy control algorithms converge in weakly communicating MDPs. As a direct extension, we show that average-reward options algorithms for temporal abstraction introduced by Wan, Naik, & Sutton (2021b) converge if the Semi-MDP induced by options is weakly communicating.

1. Introduction

In this paper, we extend the convergence results of two off-policy control algorithms for the average-reward objective from a sub-class of MDPs to the most general class of MDPs that could be solved by algorithms learning from a single stream of experience. These algorithms learn a policy that achieves the best possible average-reward rate, using data generated by some other policy that the agent may not have control of. Designing convergent off-policy algorithms for the average-reward objective is challenging. While there are several off-policy learning algorithms in the literature, the only known convergent algorithms are SSP Q-learning and RVI Q-learning, both by Abounadi, Bertsekas, & Borkar (2001), the algorithm by Ren & Krogh (2001), and Differential Q-learning by Wan, Naik, & Sutton (2021a). Others either do not have a convergence theory (Schwartz 1993, Singh 1994; Bertsekas & Tsitsiklis 1996, Das 1999) or have incorrect proof (Yang 2016, Gosavi 2004).¹

The algorithm by Ren & Krogh (2001) requires knowledge of properties of the MDP which are not typically known. The convergence of SSP Q-learning requires knowing a state that is recurrent under all policies. The convergence of the RVI Q-learning algorithm (Abounadi et al. 2001) was developed for unichain MDPs, which just means that the Markov chain induced by any stationary policy is unichain². The convergence of Differential Q-learning (Wan et al. 2021a) requires a

1. See Appendix D in Wan et al. (2021a) for a discussion about Yang’s proof and see [Appendix D](#) of this paper for a discussion about Gosavi’s proof.

2. A Markov chain is unichain if there is only one recurrent class in the Markov chain, plus a possibly empty set of transient states.

weaker assumption – the solution set of the average-reward optimality equation (formally defined later in (2)) only has one degree of freedom (all the solutions are different by a constant vector). This assumption can be satisfied if, for example, all optimal policies are unichain. It is clear that RVI Q-learning also converges under this assumption. It is not rare that the solution set of the average-reward optimality equation has more than one degree of freedom (e.g., the MDP at the bottom of Figure 1). In this case, the proofs of RVI Q-learning and Differential Q-learning would not go through. Technically, this is because both two proofs require that the uniqueness of the solution of the action values up to an additive constant (one degree of freedom) in the average-reward optimality equation, so that there is a unique equilibrium in the ordinary differential equations associated with the two algorithms.

A more general class of MDPs, called weakly communicating MDPs, may have more than a single degree of freedom in the solution set of the associated optimality equation. The definition of these MDPs is also natural: except for a possibly empty set of states that are transient under every policy, all states are reachable from every other state in a finite number of steps with a non-zero probability. It has been observed that the set of weakly communicating MDPs is the most general set of MDPs such that there exists a learning algorithm that can, using a single stream of experience, guarantee to identify a policy that achieves the optimal average reward rate in the MDP (Barlett & Tewari 2009).

In this paper, we show the convergence of RVI Q-learning and Differential Q-learning in weakly communicating MDPs, without requiring any additional assumptions compared with their original convergence theories. Two key steps in our proof are 1) showing that the solution sets of the two algorithms are non-empty, closed, bounded, and connected, and 2) showing that 0 is the unique solution for the average-reward optimality equation when all rewards are 0. With these two results, we use asynchronous stochastic approximation results by Borkar (2009) to show convergence to the solution

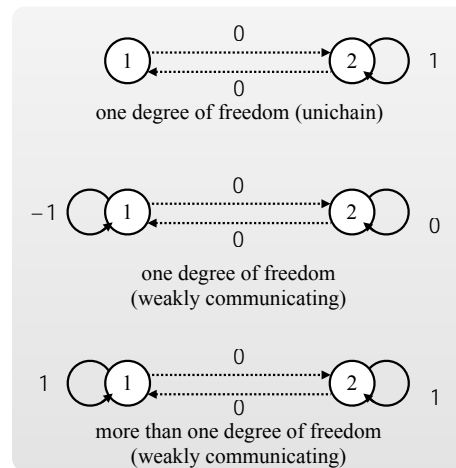


Figure 1: Examples of three different types of MDPs. In each of the three MDPs, there are two states marked by two circles respectively. There are two actions solid and dashed, both causing deterministic effects. *Top*: The solution set of q in the average-reward optimality equation ((2)) is $f q(1; \text{dashed}) = c - 1; q(2; \text{solid}) = c; q(2; \text{dashed}) = c - 2; 8c - 2 Rg$ and has one degree of freedom. The Markov chain under every stationary policy is unichain. *Middle*: The solution set of q is $f q(1; \text{solid}) = c - 3; q(1; \text{dashed}) = c - 1; q(2; \text{solid}) = c; q(2; \text{dashed}) = c - 2; 8c - 2 Rg$ and has one degree of freedom. The MDP is weakly communicating. *Down*: The solution set of q is has more than one degree of freedom. Note that both $f q(1; \text{solid}) = 0; q(1; \text{dashed}) = 2; q(2; \text{solid}) = 1; q(2; \text{dashed}) = 1g$ and $f q(1; \text{solid}) = 0; q(1; \text{dashed}) = 1; q(2; \text{solid}) = 0; q(2; \text{dashed}) = 1g$ are solutions of q in (2) and these two solutions are not different by a constant vector. The MDP is weakly communicating.

sets. As a direct extension of the above results, we also show the convergence of two algorithms that extend the Differential Q-learning algorithm to the options framework, introduced by Wan et al. (2021b), if the Semi-MDP induced by a given MDP and a given set of options is weakly communicating.

2. Preliminaries

Consider a finite Markov decision process, defined by the tuple $\mathcal{M} \doteq (S; A; R; p)$, where S is a set of states, A is a set of actions, R is a set of rewards, and $p: S \times R \times S \times A \rightarrow [0; 1]$ is the dynamics of the MDP. At each time step t , the agent observes the state of the MDP $S_t \in S$ and chooses an action $A_t \in A$ using some policy $b: A \times S \rightarrow [0; 1]$, then receives from the environment a reward $R_{t+1} \in R$ and the next state $S_{t+1} \in S$, and so on. The transition dynamics are defined as $p(s^{\prime}; r \mid s; a) \doteq \Pr(S_{t+1} = s^{\prime}; R_{t+1} = r \mid S_t = s; A_t = a)$ for all $s; s^{\prime} \in S; a \in A$, and $r \in R$. Denote the set of stationary Markov policies Π .

The reward rate of a policy π starting from a given start state S can be defined as:

$$r(\pi; s) \doteq \lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E}[R_t \mid S_0 = s; A_{0:t-1} = \pi] \quad (1)$$

Given an arbitrary MDP, the agent may not even be able to visit all states and would therefore miss the chance of learning, for every state S , a policy that achieves the optimal reward rate $\sup_{\pi} r(\pi; s)$ and the agent can at most learn an optimal policy for a set of states, each of which is reachable from every other state. Such a set of states is often called *communicating*. Formally speaking, we say a set of states *communicating*, if there exists a policy such that moving from either one state in the set to the other one in the set in a finite number of steps has a positive probability. If the entire state space of an MDP is communicating, we say the MDP *communicating*. *Weakly communicating* MDPs generalize over communicating MDPs. In weakly communicating MDPs, in addition to a *closed* communicating set of states, there is a possibly empty set of states that are transient under every policy.

For weakly communicating MDPs, there exists a unique optimal reward rate r^* , which does not depend on the start state. We say a policy is optimal if it achieves r^* regardless of the start state. The goal of an off-policy control algorithm is to learn an optimal policy from the stream of experience $\{S_t; A_t; R_{t+1}; S_{t+1}\}_{t=0}^{\infty}$ generated by a behavior policy that is not necessarily the same as the agent’s learned policy. Both RVI Q-learning and Differential Q-learning achieve this goal by solving r and q in the optimality equation:

$$q(s; a) = \sum_{s^{\prime}; r} p(s^{\prime}; r \mid s; a) (r + r^* + \max_{a^{\prime}} q(s^{\prime}; a^{\prime})); \quad \forall s \in S; a \in A \quad (2)$$

It is known that r^* is the unique solution of r and any greedy policy w.r.t. any solution of q is an optimal policy. In addition, shifting any solution of q by any constant vector results in the other solution of q . Finally, unlike in unichain MDPs, where all solutions of q are different by some constant vector, in weakly communicating MDPs, solutions of q may have multiple degrees of freedom. That is, if $q_1; q_2$ are both solutions of q , it is possible that $q_1 \in q_2 + c\mathbf{e}$; $\forall c \in \mathbb{R}$, where \mathbf{e} denotes the all-one vector.

3. Convergence Results

In this section, we present convergence theories of Differential Q-learning and RVI Q-learning in weakly communicating MDPs, and theories of the two option extensions of Differential Q-learning in weakly communicating SMDPs. Empirical results verifying the convergence of the two MDP algorithms are presented in [Appendix C](#).

Differential Q-learning updates a table of estimates $Q_t : S \times A \rightarrow \mathbb{R}$ as follows:

$$\begin{aligned} Q_{t+1}(S_t; A_t) &\stackrel{\dot{=}}{=} Q_t(S_t; A_t) + \alpha_{(t; S_t; A_t)} \delta_t \\ Q_{t+1}(s; a) &\stackrel{\dot{=}}{=} Q_t(s; a); \quad \forall s; a \notin S_t; A_t \end{aligned} \quad (3)$$

where $\alpha_{(t; S_t; A_t)}$ is the number of times $S_t; A_t$ has been visited before time step t , $\alpha_{(t; S_t; A_t)}^{-1}$ is a step-size sequence, and δ_t , the temporal-difference (TD) error, is:

$$\delta_t \stackrel{\dot{=}}{=} R_{t+1} - R_t + \max_a Q_t(S_{t+1}; a) - Q_t(S_t; A_t); \quad (4)$$

where R_t is a scalar estimate of r , updated by:

$$R_{t+1} \stackrel{\dot{=}}{=} R_t + \alpha_{(t; S_t; A_t)} \delta_t \quad (5)$$

and α is a positive constant.

We now present the convergence theory of Differential Q-learning. We first state the required assumptions, which are also required by the original convergence theory of Differential Q-learning by Wan et al. (2021a).

Assumption 1 For all $n \geq 0$, $\alpha_n > 0$, $\sum_{n=0}^{\infty} \alpha_n = 1$, and $\sum_{n=0}^{\infty} \alpha_n^2 < 1$.

Assumption 2 Let $\lfloor \cdot \rfloor$ denote the integer part of (\cdot) , for $x \in (0; 1)$, $\sup_n \frac{\lfloor xn \rfloor}{n} < 1$ and $\frac{\sum_{n=0}^{\lfloor ym \rfloor} \alpha_n}{\sum_{n=0}^m \alpha_n} \rightarrow 1$ uniformly in $y \in [x; 1]$.

Assumption 3 There exists $\epsilon > 0$ such that

$$\liminf_{n \rightarrow \infty} \frac{\alpha_n(s; a)}{\alpha_{n+1}(s; a)} > \epsilon;$$

a.s., for all $s \in S; a \in A$. Furthermore, for all $x > 0$, let $N(n; x) = \min \{m > n : \sum_{k=n+1}^m \alpha_k < x\}$; the limit $\lim_{n \rightarrow \infty} \frac{\sum_{k=N(n; x); s; a}^{\infty} \alpha_k}{\sum_{k=N(n; x); s^0; a^0}^{\infty} \alpha_k} = 1$ exists a.s. for all $s; s^0 \in S; a; a^0 \in A$.

Theorem 1 If \mathcal{M} is communicating and Assumptions 1–3 hold, the Differential Q-learning algorithm (Equations 3–5) converges, almost surely, R_t to r , Q_t to the set of solutions of (2) and

$$r = R_0 = \prod_{s; a} q(s; a) \prod_{s; a} Q_0(s; a); \quad (6)$$

and $r(\pi_t; s)$ to r , for all $s \in S$, where π_t is any greedy policy w.r.t. Q_t .

Remark: If the MDP is weakly communicating, that is, it contains transient states, the agent eventually reaches the closed communicating state and never returns to the transient states. Elements in Q_η that are associated with the closed communicating set converge to a set that depends on the values of Q and R when the MDP reaches the closed communicating set for the first time. Other elements in Q_η would only be visited for a finite number of times and can not be guaranteed to converge to their correct values by any learning algorithm. Other conclusions of the theorem remain unchanged. This observation on weakly communicating MDPs also applies to Theorems 2–4.

The update rules of RVI Q-learning are

$$\begin{aligned} Q_{t+1}(S_t; A_t) &\stackrel{\dot{=}}{=} Q_t(S_t; A_t) + \alpha_{(t; S_t; A_t)} (r_t(S_t; A_t) - Q_t(S_t; A_t)); \\ Q_{t+1}(s; a) &\stackrel{\dot{=}}{=} Q_t(s; a); \quad \forall s; a \notin S_t; A_t; \end{aligned} \quad (7)$$

where

$$\alpha_{(t; S_t; A_t)} \stackrel{\dot{=}}{=} R_{t+1} - f(Q_t) + \max_a Q_t(S_{t+1}; a) - Q_t(S_t; A_t); \quad (8)$$

and $f : S \times A \rightarrow \mathbb{R}$ satisfies the following assumption.

Assumption 4 1) f is L -Lipschitz, 2) there exists a positive scalar u s.t. $f(e) = u$ and $f(x + ce) = f(x) + cu$, and 3) $f(cx) = cf(x); \forall c \geq 0$.

Theorem 2 If \mathcal{M} is communicating and Assumptions 1–4 hold, then the RVI Q-learning algorithm (Equations 7–8) converges, almost surely, R_t to r , Q_t to the set of solutions of (2) and

$$r = f(q); \quad (9)$$

and $r(s) = f(q)$ for all $s \in S$, where q is any greedy policy w.r.t. Q_t .

4. Conclusions

In this paper, we provide, for the first time, convergence results of off-policy average-reward control algorithms in weakly communicating MDPs, which are known to be the most general class of MDPs in which it is possible that a learning algorithm can guarantee to obtain an optimal policy. Specifically, we show two existing algorithms, RVI Q-learning and Differential Q-learning, converge in weakly communicating MDPs. As an extension, in Appendix A, we also showed two off-policy average-reward options learning algorithms converge if the SMDP induced by the options is weakly communicating.

Acknowledgements

The authors were generously supported by DeepMind, Amii, NSERC, and CIFAR. The authors wish to thank Huizhen Yu and Abhishek Naik for discussing several important related papers and discussing ideas to address the technical challenges in the proof. Computing resources were provided by Compute Canada.

References

- Abounadi, J., Bertsekas, D., Borkar, V. S. (2001). Learning Algorithms for Markov Decision Processes with Average Cost. *SIAM Journal on Control and Optimization*.
- Bartlett, P., & Tewari, A. (2009). REGAL: a regularization based algorithm for reinforcement learning in weakly communicating MDPs. *Uncertainty in Artificial Intelligence*.
- Bertsekas, D. P. (2007). *Dynamic Programming and Optimal Control third edition, volume II*. Athena Scientific.
- Bertsekas, D. P., Tsitsiklis, J. N. (1996). *Neuro-dynamic Programming*. Athena Scientific.
- Borkar, V. S. (1998). Asynchronous Stochastic Approximations. *SIAM Journal on Control and Optimization*.
- Borkar, V. S. (2009). *Stochastic Approximation: A Dynamical Systems Viewpoint*. Springer.
- Das, T. K., Gosavi, A., Mahadevan, S. Marchallick, N. (1999). Solving semi-Markov decision problems using average reward reinforcement learning. *Management Science*.
- Gosavi, A. (2004). Reinforcement learning for long-run average cost. *European Journal of Operational Research*.
- Puterman, M. L. (1994). *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons.
- Ren, Z., Krogh, B. H. (2001). Adaptive control of Markov chains with average cost. *IEEE Transactions on Automatic Control*.
- Schwartz, A. (1993). A reinforcement learning method for maximizing undiscounted rewards. In *Proceedings of the International Conference on Machine Learning*.
- Schweitzer, P. J. (1971). Iterative solution of the functional equations of undiscounted Markov renewal programming. *Journal of Mathematical Analysis and Applications*.
- Schweitzer, P. J., & Federgruen, A. (1978). The Functional Equations of Undiscounted Markov Renewal Programming. *Mathematics of Operations Research*.
- Singh, S. P. (1994). Reinforcement learning algorithms for average-payoff Markovian decision processes. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Sutton, R. S., Precup, D., Singh, S. (1999). Between MDPs and Semi-MDPs: A Framework for Temporal Abstraction in Reinforcement Learning. *Artificial Intelligence*.
- Sutton, R. S., Barto, A. G. (2018). *Reinforcement Learning: An Introduction*. MIT Press.
- Wan, Y., Naik, A., Sutton, R. S. (2021a). Learning and Planning in Average-Reward Markov Decision Processes. *International Conference on Machine Learning*.
- Wan, Y., Naik, A., Sutton, R. S. (2021b). Average-Reward Learning and Planning with Options. *Conference on Neural Information Processing Systems*.

Appendix A. Learning with Options

If the agent has a set of options, it may choose to execute these options. Each option o in O has two components: the option's *policy* $o : A \times S \rightarrow [0;1]$, and the termination probability $o : S \rightarrow [0;1]$. For simplicity, for any $s \in S; o \in O$, we use $(a | j s; o)$ to denote $o(a; s)$ and $(s; o)$ to denote $o(s)$. If the agent executes option o at state s , the option's policy is followed, until the option terminates. Let L be the set of all possible lengths of options and \hat{R} be the set of all possible cumulative rewards. Note that L and \hat{R} are possibly countably infinite. Let $\hat{p}(s^d; r; l | j s; o)$ be, when executing option o starting from state s , the probability of terminating at state s^d , with cumulative reward r and length l . Formally, for any $s; s^d \in S; o \in O; r \in \hat{R}; l \in L$, \hat{p} can be defined recursively in the following way:

$$\hat{p}(s^d; r; l | j s; o) \doteq \prod_a (a | j s; o) \prod_{s':F} p(s; F | j s; a) \\ [(s; o) \mathbf{1}(s = s^d; F = r; l = 1) + (1 - (s; o)) \hat{p}(s^d; r - F; l - 1 | j s; o)]; \quad (10)$$

where $\mathbf{1}$ is an indicator function.

An MDP M and a set of options O results in a Semi-MDP (SMDP) $\hat{M} \doteq (S; O; L; \hat{R}; \hat{p})$.

Given an MDP and a set of options, if the agent chooses options using a meta policy, which is a policy that chooses from options, $\pi : O \times S \rightarrow [0;1]$ and executes these options, we denote the sequence of option transitions by $\pi; \hat{S}_0; \hat{O}_0; \hat{R}_{0+1}; \hat{S}_{0+1}; \dots$. For the associated SMDP, the *reward rate* of π given a start state s can be defined as $r^C(\pi; s) \doteq \lim_{t \rightarrow \infty} \frac{1}{t} E[\sum_{i=1}^t R_i | j S_0 = s]$ or $r(\pi; s) \doteq \lim_{n \rightarrow \infty} \frac{1}{n} E[\sum_{i=0}^n \hat{R}_i | j \hat{S}_0 = s] = E[\sum_{i=0}^n \hat{L}_i | j \hat{S}_0 = s]$. Both limits exist and are equivalent (by Puterman's (1994) Propositions 11.4.1 and 11.4.7) under the following assumption:

Assumption 5 *For each option $o \in O$, when executing the option, there is a non-zero probability of terminating the option after at most $|S|$ stages, regardless of the state at which this option is initiated.*

Proposition 1 *Under Assumption 5, the expected value as well as the variance of the execution time and cumulative reward of every option at each state exist and are finite.*

We say an SMDP is weakly communicating if the MDP with state space S , action space O , reward space \hat{R} , and transition function $\hat{p}(s; r; l | j s; o)$ is weakly communicating. Just as in the MDP setting, if the SMDP is weakly communicating, the optimal reward rate $\hat{r} \doteq \sup_{\pi \in \hat{\Pi}} r(\pi; s)$, where $\hat{\Pi}$ denotes the set of stationary Markov meta policies, does not depend on the start state s . In addition, the solutions of q may not be different by a constant vector. Given an MDP and a set of options, the goal of the off-policy control problem is to find a policy that achieves \hat{r} . Inter-option Differential Q-learning achieves this goal by solving the *optimality* equation for SMDPs (Puterman 1994):

$$q(s; o) = \sum_{s^d; r; l} \hat{p}(s^d; r; l | j s; o) (r - r + l + \max_{o^d} q(s^d; o^d)); \quad (11)$$

where q and r denote estimates of the option-value function and the reward rate respectively. Just as in the MDP setting, r has \hat{r} as its unique solution, and solutions of q may not be different by a constant vector.

Intra-option Differential-learning finds an optimal policy by solving the *intra-option optimality* equation.

$$q(s; o) = \underset{a}{\times} (a; s; o) \underset{s^o; r}{\times} p(s^o; r; j; s; a) (r + u^q(s^o; o)); \quad \forall s \in S; o \in O; \quad (12)$$

where

$$u^q(s^o; o) \doteq 1 - (s^o; o) q(s^o; o) + (s^o; o) \max_{o'} q(s^o; o'); \quad (13)$$

The following proposition shows that the set of solutions of (12) is the same as that of (11).

Proposition 2 *Any solution of (11) is also a solution of (12) and vice versa.*

Now consider option extensions of Differential Q-learning. Given an SMDP $\hat{\mathcal{M}} = (S; O; L; \hat{R}; \hat{p})$, inter-option Differential Q-learning maintains estimates of option values, and, inspired by Schweitzer (1971), updates estimates using scaled TD errors:

$$Q_{n+1}(\hat{S}_n; \hat{O}_n) \doteq Q_n(\hat{S}_n; \hat{O}_n) + \frac{1}{(n; \hat{S}_n; \hat{O}_n)} (R_n - L_n(\hat{S}_n; \hat{O}_n) Q_n(\hat{S}_n; \hat{O}_n)); \quad (14)$$

$$Q_{n+1}(s; o) \doteq Q_n(s; o); \quad \forall s; o \notin \hat{S}_n; \hat{O}_n; \quad (15)$$

$$R_{n+1} \doteq R_n + \frac{1}{(n; \hat{S}_n; \hat{O}_n)} (R_n - L_n(\hat{S}_n; \hat{O}_n) R_n);$$

where $(n; \hat{S}_n; \hat{O}_n)$ is the number of visits to state-option pair $(\hat{S}_n; \hat{O}_n)$ before stage n , $L_n(\cdot)$ comes from an additional vector of estimates $L: S \times O \rightarrow \mathbb{R}$ that approximate the expected lengths of state-option pairs, updated by:

$$L_{n+1}(\hat{S}_n; \hat{O}_n) \doteq L_n(\hat{S}_n; \hat{O}_n) + \frac{1}{(n; \hat{S}_n; \hat{O}_n)} (L_n - L_n(\hat{S}_n; \hat{O}_n)); \quad (16)$$

where $f_n g$ is the other step-size sequence. The TD-error δ_n in (14) and (15) is

$$\delta_n \doteq \hat{R}_n - L_n(\hat{S}_n; \hat{O}_n) R_n + \max_{o'} Q_n(\hat{S}_{n+1}; o) - Q_n(\hat{S}_n; \hat{O}_n); \quad (17)$$

Theorem 3 *If $\hat{\mathcal{M}}$ is communicating, Assumptions 1–3 and 5 hold, except for using $(n; s; o)$ instead of $(t; s; a)$, and that $0 < \beta_n < 1$, $\beta_n \beta_{n-1} > 1$, and $\beta_n^2 < 1$, inter-option Differential Q-learning (Equations 14–17) converges, almost surely, Q_n to the set of solutions of (11) and*

$$\hat{r} = R_0 = \underset{s; o}{\times} q(s; o) \underset{s; o}{\times} Q_0(s; o); \quad (18)$$

R_n to \hat{r} , and $r(n; s)$ to \hat{r} where π_n is a greedy policy w.r.t. Q_n .

Intra-option Differential Q-learning also maintains estimates of option values. However, instead of updating the estimates using option transitions, it updates for all options using each action transition $(S_t; O_t; A_t; R_{t+1}; S_{t+1})$.

$$Q_{t+1}(S_t; o) \doteq Q_t(S_t; o) + \alpha \frac{A_t(S_t; o)}{A_t(S_t; O_t)} (R_{t+1} - Q_t(S_t; o)); \quad \forall o \in O; \quad (19)$$

$$Q_{t+1}(s; o) \doteq Q_t(s; o); \quad \forall s \in S; o \in O; \quad (20)$$

$$R_{t+1} \doteq R_t + \alpha \sum_{o \in O} \frac{A_t(S_t; o)}{A_t(S_t; O_t)} (R_{t+1} - Q_t(S_t; o));$$

where α is a step-size sequence, $\frac{A_t(S_t; o)}{A_t(S_t; O_t)}$ is the importance sampling ratio, and:

$$u^{Q_t}(S_{t+1}; o) \doteq R_{t+1} - R_t + u^{Q_t}(S_{t+1}; o) - Q_t(S_t; o); \quad (21)$$

where u^{Q_t} is defined in (13).

Theorem 4 *If \mathcal{M} is communicating, Assumptions 1–3 and 5 hold, except for using $(t; s; o)$ instead of $(t; s; a)$, intra-option Differential Q-learning (Equations 19–21) converges, almost surely, Q_t to the set of solutions of (11) and (18), R_t to \hat{r} , and $r(t; s)$ to \hat{r} where t is a greedy policy w.r.t. Q_t .*

Appendix B. Proofs

B.1. Proof of Proposition 1

Note that the execution time of each option $o \in O$ is the return of executing this option's policy in a stochastic shortest path MDP (SSP-MDP, Bertsekas 2007) with state space $S + f^?g$ ($?$ is the ‘‘terminal’’ state of the SSP-MDP), action space A and transition function \hat{p} satisfying:

$$\hat{p}(s^o; 1 | j; s; a) \doteq (1 - \sum_{o' \in O} p(s^o; r | j; s; a)) p(s^o; r | j; s; a); \quad \forall s; s^o \notin ?; a \in A;$$

$$\hat{p}(?; 1 | j; s; a) \doteq \sum_{o \in O} p(s^o; r | j; s; a); \quad \forall s \notin ?; a \in A;$$

$$\hat{p}(?; 0 | j; ?; a) \doteq 1; \quad \forall a \in A;$$

Also, note that by Assumption 5, the option's policy is a ‘proper’ policy (Bertsekas 2007) in the SSP-MDP. That is, when using the policy, the MDP reaches the terminal state eventually regardless of the start state. Because the expected value of every proper policy of an SSP-MDP exists and is finite (Section 2.1 of Bertsekas (2007)), the expected value of the execution time of option o exists. The existence of the variance can be shown using similar arguments as those used to show the existence of the expectation in Section 2.1 of Bertsekas (2007).

Similarly, the cumulative reward of each option o is the return of executing o 's policy in an SSP-MDP with state space $S + f^?g$, action space A , and transition function \hat{p} satisfying:

$$\hat{p}(s^o; r | j; s; a) \doteq (1 - \sum_{o' \in O} p(s^o; r | j; s; a)) p(s^o; r | j; s; a); \quad \forall s; s^o \notin ?; a \in A$$

$$\hat{p}(?; r | j; s; a) \doteq \sum_{o \in O} p(s^o; r | j; s; a); \quad \forall s \notin ?; a \in A$$

$$\hat{p}(?; 0 | j; ?; a) \doteq 1; \quad \forall a \in A;$$

Again the option's policy is proper and the expected value of the cumulative reward of option o exists. So does the variance of the cumulative reward.

B.2. Proof of Proposition 2

By the definition of \hat{p} of the SMDP induced by choosing options in an MDP,

$$\hat{p}(s; F; t; j; s; o) = \prod_a (a; j; s; o) \prod_{s^d; r} p(s^d; r; j; s; a)$$

$$[(s^d; o) \mathbf{I}(s = s^d; F = r; t = 1) + (1 - (s^d; o)) \hat{p}(s; F; r; t - 1; j; s^d; o)]$$

$$\begin{aligned} & q(s; o) \\ &= \prod_{s; F; t} \hat{p}(s; F; t; j; s; o) (F - tr + \max_{o'} q(s^d; d^d)) \\ &= \prod_{s; F; t} \prod_a (a; j; s; o) \prod_{s^d; r} p(s^d; r; j; s; a) \\ & [(s^d; o) \mathbf{I}(s = s^d; F = r; t = 1) + (1 - (s^d; o)) \hat{p}(s; F; r; t - 1; j; s^d; o)] (F - tr + \max_{o'} q(s^d; d^d)) \\ &= \prod_a (a; j; s; o) \prod_{s^d; r} p(s^d; r; j; s; a) (s^d; o) (r - r + \max_{o'} q(s^d; d^d)) \\ &+ \prod_a (a; j; s; o) \prod_{s^d; r} p(s^d; r; j; s; a) (1 - (s^d; o)) \prod_{s; F; t} \hat{p}(s; F; r; t - 1; j; s^d; o) (F - tr + \max_{o'} q(s^d; d^d)) \\ &= \prod_a (a; j; s; o) \prod_{s^d; r} p(s^d; r; j; s; a) (s^d; o) (r - r + \max_{o'} q(s^d; d^d)) \\ &+ \prod_a (a; j; s; o) \prod_{s^d; r} p(s^d; r; j; s; a) (1 - (s^d; o)) \prod_{s; F; t} \hat{p}(s; F; t; j; s^d; o) (F + r - (t + 1)r + \max_{o'} q(s^d; d^d)) \\ &= \prod_a (a; j; s; o) \prod_{s^d; r} p(s^d; r; j; s; a) (s^d; o) (r - r + \max_{o'} q(s^d; d^d)) \\ &+ \prod_a (a; j; s; o) \prod_{s^d; r} p(s^d; r; j; s; a) (1 - (s^d; o)) @_r r + \prod_{s; F; t} \hat{p}(s; F; t; j; s^d; o) (F - tr + \max_{o'} q(s^d; d^d)) \text{A} \\ &= \prod_a (a; j; s; o) \prod_{s^d; r} p(s^d; r; j; s; a) (s^d; o) (r - r + \max_{o'} q(s^d; d^d)) \\ &+ \prod_a (a; j; s; o) \prod_{s^d; r} p(s^d; r; j; s; a) (1 - (s^d; o)) r - r + q(s^d; o) \\ &= \prod_a (a; j; s; o) \prod_{s^d; r} p(s^d; r; j; s; a) (r - r + u^{\hat{p}}(s^d; o)) \end{aligned}$$

Therefore any solution of (11) must be a solution of (12)-(13) and vice versa.

B.3. Proof sketch of Theorem 1-Theorem 4

In this section, we sketch the proof of Theorems 1-4. It has been shown that all four algorithms introduced above are special cases of the *General RVI Q algorithm* (Wan et al. 2021a,b). They also showed that General RVI Q converges under an assumption that is not satisfied for weakly communicating MDPs/SMDPs. In order to show convergence for weakly communicating MDPs/SMDPs, we replace this assumption with three weaker assumptions that are satisfied for these MDPs/SMDPs. All other assumptions are the same as those used by Wan et al. (2021a,b) and can be verified for all four algorithms using their arguments. We present General RVI Q and prove its convergence with the three new assumptions in Appendix B.8. The next step of the proof would be verifying the three new assumptions when casting General RVI Q to each of the four algorithms. This should be straightforward given our Theorem 5 and Lemma 4. We defer this part to Appendix B.9. Given that the three assumptions are verified, we have the conclusion part of the convergence theorem of General RVI Q holds for each of the four algorithms. The convergence of the reward rates of greedy policies w.r.t. the action/option-values follows the convergence of these values and is shown in Appendix B.10.

B.4. Characterization of the Solution Set

In this section, we characterize the sets that the algorithms described in the previous section converge to. This section plays a key role in showing their convergence.

We consider the set of solutions of q in the SMDP optimality equation ((11)) and

$$\hat{r} = f(q); \tag{22}$$

where $f : S \times O \rightarrow \mathbb{R}$ satisfies Assumption 4. It is clear that (11) generalizes over (2) and (22) generalizes over (6), (9), and (18). And thus the characterization of Q_1 applies to the sets that action/option values in the aforementioned algorithms are claimed to converge to in Theorems 1-4.

It is known that if the SMDP is weakly communicating, (11) has \hat{r} as its unique solution of r . For q , it has been shown by Schweitzer & Federgruen (1978) (we will refer to this work multiple times and thus we use a shorthand ‘‘S&F’’ for simplicity from now on) in their Theorem 4.2 that the set of solutions of q in (11) is closed, unbounded, connected, and possibly non-convex. The next theorem characterizes Q_1 .

Theorem 5 *If the SMDP is weakly communicating and Assumption 4 holds, Q_1 is non-empty, closed, bounded, connected, and possibly non-convex.*

Proof

First, Q_1 is non-empty. To see this point, note that for any solution of q in (11), $q, q + ce$ is also a solution for any $c \geq 0$ and thus there must be a c such that (22) holds because $f(x + ce) = f(x) + cu$ for any x, c .

Q_1 is closed because the set of solutions of q in (11) is closed by S&F, the set of solutions of q in (22) is closed because f is Lipschitz and is thus continuous, and the intersection of two closed sets is closed.

Boundedness

We now show that Q_1 is bounded. For any $q \in Q_1$, let $v(s) \doteq \max_o q(s; o)$. Rewrite the option-value optimality equation (Equation (11)) using v instead of q , we have,

$$v(s) = \max_o \sum_{s^j; r; l} p(s^j; r; l | s; o) (r + v(s^j)); \quad \forall s \in S. \quad (23)$$

The above equation is known as the state-value optimality equation.

It is easy to verify that

$$q(s; o) = \sum_{s^j; r; l} p(s^j; r; l | s; o) (r + v(s^j)); \quad (24)$$

Using this fact, rewrite Equation (22) using v instead of q , we have,

$$f^* = f(F + P^*v); \quad (25)$$

where

$$F(s; o) \doteq \sum_{s^j; r; l} p(s^j; r; l | s; o) (r + f^*);$$

$$P(s; o; s^j) \doteq \sum_{r; l} p(s^j; r; l | s; o);$$

Denote the set of solutions of v in (23) by V . Denote the set of solutions of v in Equation (23) and Equation (25) by V_1 . If V_1 is bounded, Q_1 is also bounded because any $q \in Q_1$ can be obtained from a solution of $v \in V_1$ with a linear operation in view of Equation (24).

In order to show boundedness, We will need the following two lemmas to proceed. These two lemmas are similar to Theorem 4.1 (c) and Theorem 5.1 in S&F, except that 1) the 'max' operates over the set of all optimal policies, instead of the set of all deterministic optimal policies as in Theorem 4.1 (c), and 2) we consider the set of weakly communicating SMDPs while S&F considers general multi-chain SMDPs. The proofs are also essentially the same. For completeness, we provide the proofs for these two lemmas in Sections B.5, B.6.

To formally state the Lemma 1, we will first introduce some definitions.

For any $\pi \in \Pi$, let P^π denote the $(S \times S) \times (S \times S)$ transition probability matrix under policy π . That is,

$$P^\pi(s; s^j) \doteq \sum_{o; r; l} p(s^j; r; l | s; o) \pi(o | s); \quad (26)$$

Let P^{π^*} be the *limiting matrix* of P^π , which is the Cesaro limit of the sequence $f P^{\pi^*} g_{i=1}^1$:

$$P^{\pi^*} \doteq \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=0}^{n-1} P^{\pi^* i}; \quad (27)$$

Because S is finite, the Cesaro limit exists and P^{π^*} is a stochastic matrix (has row sums equal to 1).

Let $I(s) \doteq \sum_{o; s^j; r; l} p(s^j; r; l | s; o) l$. And let the fundamental matrix $Z \doteq (I - P^{\pi^*} + P^{\pi^*})^{-1} = I + \lim_{n \rightarrow \infty} \sum_{i=1}^n (P^{\pi^* i} - P^{\pi^*})$.

Lemma 1 *If the SMDP is weakly communicating, v is a solution of (23) if and only if*

$$v(s) = \max_{\pi \in \hat{\Pi}} [r + \int P^\pi v](s); \forall s \in S: \tag{28}$$

In order to state Lemma 2 formally, we need the following definitions. Define the Bellman error for a state $s \in S$ given a policy $\pi \in \hat{\Pi}$ and some $v \in \mathbb{R}^{|S|}$, $b_{v,\pi}(s)$, as follows:

$$b_{v,\pi}(s) \doteq [r + \int P^\pi v - v](s):$$

Define R as the set of recurrent states for P^π . That is,

$$R \doteq \{s \in S \mid \exists \ell \geq 1, P^\pi(s, s) > 0\}:$$

Let $\hat{\Pi}$ denote the set of optimal meta policies. Define R as the set of states that are recurrent under some optimal meta policy:

$$R \doteq \{s \in S \mid \exists \pi \in \hat{\Pi}, s \in R_\pi\}:$$

By Theorem 3.2 (b) in S&F, there exists a policy $\pi \in \hat{\Pi}$ such that $R = R_\pi$. For any $\pi \in \hat{\Pi}$, let $n(\pi)$ be the number of recurrent classes for P^π . Further, define the least number of recurrent classes induced by any optimal meta policy that induces the set of recurrent states R :

$$n \doteq \min_{\pi \in \hat{\Pi}} n(\pi); R = R_\pi:$$

Denote $\hat{\Pi}_n$ as the set of policies that have R as their sets of recurrent states and that have n recurrent classes. That is,

$$\hat{\Pi}_n \doteq \{\pi \in \hat{\Pi} \mid R = R_\pi; n(\pi) = n\}:$$

Theorem 3.2 (d) by S&F shows that all policies within $\hat{\Pi}_n$ share the same collection of recurrent classes. Denote the collection of recurrent classes as $\{R^1, \dots, R^n\}$. The following lemma shows that the solution set of (23) has n degrees of freedom.

Lemma 2 *If the SMDP is weakly communicating, suppose v and $v + x$ are both solutions of (23), then there exists n constants y_1, y_2, \dots, y_n such that*

$$x(s) = y_i; i \in \{1, \dots, n\}; s \in R^i \tag{29}$$

$$x(s) = \max_{\pi \in \hat{\Pi}} [Z_{b_{v,\pi}}](s) + \sum_{i=1}^n \int_{s' \in R^i} P^\pi(s, s') A y_i; s \in S \setminus R; \tag{30}$$

$$y_i = [Z_{b_{v,\pi}}](s) + \sum_{i=1}^n \int_{s' \in R^i} P^\pi(s, s') A y_i; i = 1, \dots, n; s \in R^i; \pi \in \hat{\Pi}_n; \tag{31}$$

For any $\epsilon \geq \epsilon_1/2$; $\delta \in [0, 1]$, note that there exists a policy π such that R is the only one recurrent class under π . To see this point, note that the SMDP is weakly communicating, and thus we can modify π to obtain a new meta policy such that all states except for those in R are transient.

Based on the above observation and Lemma 2, for any $v \in V$, we have for any given $\epsilon \geq \epsilon_1/2$; $\delta \in [0, 1]$; $\gamma \in [0, 1]$, there exists a π , such that

$$\begin{aligned} y &= \max_{s \in R} Z(\pi) b_{v; \pi}(s) + \sum_{s' \in R} P^{\pi}(s; s') y \\ &= \max_{s \in R} Z(\pi) b_{v; \pi}(s) + y; \quad \delta = 1; \quad \gamma \in [0, 1] \end{aligned}$$

The first term $\max_{s \in R} Z(\pi) b_{v; \pi}(s)$ is a constant given v and π . Therefore we see that, for any other solution $v + x$ of (23), if y is arbitrarily large then $y; \delta = 1; \gamma \in [0, 1]$ should also be arbitrarily large. This would violate the Lipschitz assumption on f . To see this point, let

$$f(v) \doteq f(F + P v); \quad (32)$$

Let L be a Lipschitz constant of f . L is also a Lipschitz constant of f because P is a stochastic matrix and is thus a non-expansion. Choose a $v_1 \in V_1$ and a $v_2 \in V_1$. a $v_2 \in V_1$, denote $m = kv_1 - v_2k$. Choose a sufficiently large $c > 0$ such that $cu > Lkv_1 + ce - v_2k = Lkv_1 + ce - v_2 - cek = Lm$, where $v_2 \doteq v_2 + ce$. Given this choice of c , using $f(v_1) = f(v_2) = \hat{r}$ and $f(v_1 + ce) = f(v_1) + cu$, we have $f(v_1 + ce) - f(v_2) = f(v_1) + cu - f(v_2) = cu > Lkv_1 + ce - v_2k$. This inequality suggests that f is not Lipschitz continuous with a Lipschitz constant L and thus violates our assumption. Because the choice of π is arbitrary, V_1 is upper bounded.

In addition, because the choice of π is arbitrary, we have for any $\epsilon \geq \epsilon_1/2$; $\delta \in [0, 1]$; $\gamma \in [0, 1]$,

$$y = \max_{\delta \in [0, 1]; \gamma \in [0, 1]} \max_{s \in R} Z(\pi) b_{v; \pi}(s) + y$$

If y is chosen to be arbitrarily small then y should also be arbitrarily small for all $\delta = 1; \gamma \in [0, 1]$ but again this is not allowed due to (25) for the same reason as mentioned in the previous paragraph. Therefore $y; \delta \in [0, 1]; \gamma \in [0, 1]$ can not be arbitrarily small. Thus V_1 is lower bounded. Combining the upper bound and lower bound, V_1 is bounded. Therefore Q_1 is also bounded.

Connectedness

We now show that Q_1 is connected. To this end, again it is enough to show that V_1 is connected.

Define a function that takes a $v \in V$ as input and produces an element in V_1 as output. Specifically, let $g: V \rightarrow V_1$ with $g(v) = v + xe$, where x is the solution of $f(v + xe) = \hat{r}$ and f is defined in (32). Note that x is unique given v because $\hat{r} = f(v + xe) = f(v) + xu$ and thus $x = (\hat{r} - f(v))/u$.

We now show that g is Lipschitz continuous. Consider any $v_1, v_2 \in V$. Let x_1, x_2 satisfy $v_1 + x_1 e = g(v_1)$ and $v_2 + x_2 e = g(v_2)$ respectively. Again x_1, x_2 are unique given v_1, v_2 . Note that

$$\begin{aligned} & f(v_1 + x_1 e) - f(v_2 + x_2 e) \\ &= f(v_1 + x_1 e) - f(v_2 + x_2 e) + (x_1 - x_2)u \\ &= \hat{f} - \hat{f} + (x_1 - x_2)u \\ &= (x_1 - x_2)u \end{aligned}$$

$$\begin{aligned} \|x_1 - x_2\| &= \|f(v_1 + x_1 e) - f(v_2 + x_2 e)\|/u \\ &= L \|v_1 - v_2\| \\ &= L \|v_1 - v_2\| \end{aligned}$$

$$\|g(v_1) - g(v_2)\| = \|v_1 + x_1 e - v_2 - x_2 e\| = \|v_1 - v_2 + (x_1 - x_2)e\| = (1 + L) \|v_1 - v_2\|$$

Therefore g is Lipschitz continuous with Lipschitz constant $1 + L$.

Finally, because V is connected and the image of any continuous function on a connected set is connected, $g(V)$ is connected. Note that every point in $g(V)$ belongs to V_1 by definition. Every point in V_1 also belongs to $g(V)$. To see this point, pick any $v \in V_1$, we can see that $v \in V$ and that $g(v) = v$ (note that $x = 0$ given that $v \in V_1$). Thus $v \in g(V)$. Therefore $V_1 = g(V)$ is connected.

Given that V_1 is connected, Q_1 should also be connected because Q_1 is a linear transformation of V_1 (see Equation 24).

Non-convexity

We now show that both V_1 and Q_1 are not necessarily convex. We will show this point by constructing a counter-example, which involves the communicating MDP shown in the left sub-figure of Figure 2. The optimal reward rate for the MDP is 0.

Let $f(q) = \sum_{s,a} q(s;a)$. Such a choice of f satisfies the assumption on f in Theorem 5. Then by definition, for any $v \in V_1$, we have

$$f(F + Pv) = \hat{f} :$$

For the three-state MDP considered here, $\hat{f} = 0$.

$$\begin{aligned} f(F + Pv) &= \sum_{s,a} r(s;a) + p(s^j | s;a) v(s^j) \\ &= 3 + 2v(1) + 3v(2) + v(3) \end{aligned}$$

Therefore,

$$2v(1) + 3v(2) + v(3) = 3 :$$

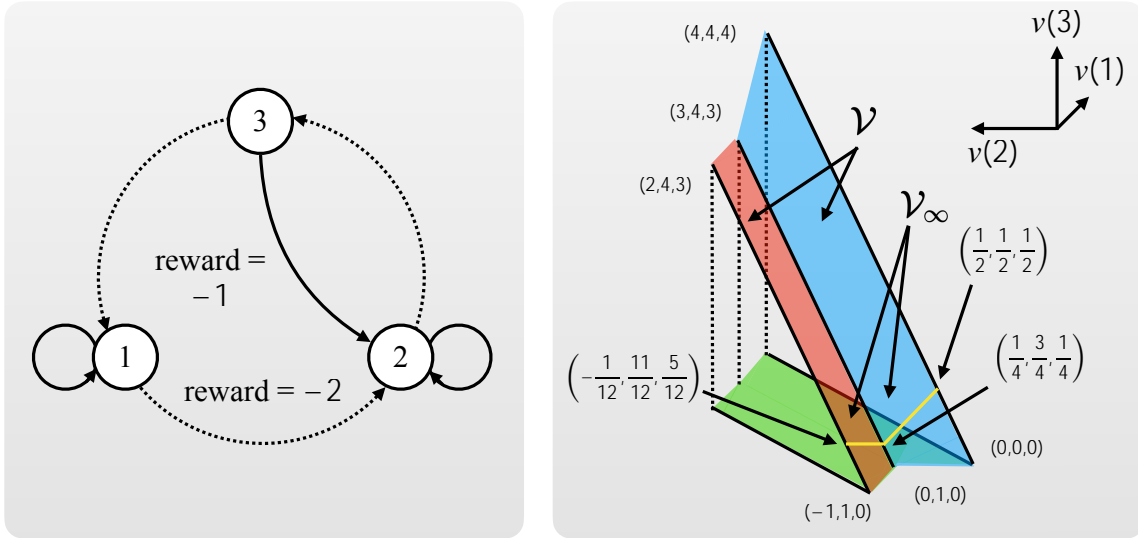


Figure 2: Illustration example. *Left*: the example MDP. There are three states marked by three circles respectively. There are two actions solid and dashed, both have deterministic effects. Taking action solid at state 3 results in a reward of -1 . Taking action dashed at state 1 results in a reward of -2 . All other rewards are 0. *Right*: a graphical explanation of V ; V_1 . The two yellow line segments together represent the solution set V_1 , which is not convex. The red and blue regions together represent V .

In addition to the above equality, V_1 needs to satisfy the state-value optimality equation (Equation (23)). Therefore for any $v \in V_1$,

$$\begin{aligned} v(1) &= \max(v(1); -2 + v(2)); \\ v(2) &= \max(v(2); v(3)); \\ v(3) &= \max(v(1); -1 + v(2)); \end{aligned}$$

which implies

$$\begin{aligned} v(1) &= -2 + v(2); \\ v(2) &= v(3); \\ v(3) &= \max(v(1); -1 + v(2)); \end{aligned}$$

Therefore

$$\begin{aligned} V_1 &= \{v \in \mathbb{R}^3 \mid v(1) = -2 + v(2); v(2) = v(3); v(3) = \max(v(1); -1 + v(2)); \\ &2v(1) + 3v(2) + v(3) = 3g\}; \end{aligned}$$

Graphically, V_1 corresponds to the two connected yellow line segments in the right sub-figure in Figure 2. From the figure, we see that V_1 is not convex.

Let s and d denote solid and dashed respectively. Consider any $q \in Q_1$, in view of (22),

$$q(1; s) + q(1; d) + q(2; s) + q(2; d) + q(3; s) + q(3; d) = 0;$$

In addition to the above equality, Q_1 needs to satisfy the action-value optimality equation (Equation (2)). Therefore for any $q \in Q_1$,

$$\begin{aligned} q(1; s) &= 0 + \max(q(1; s); q(1; d)) \\ q(1; d) &= 2 + \max(q(2; s); q(2; d)) \\ q(2; s) &= 0 + \max(q(2; s); q(2; d)) \\ q(2; d) &= 0 + \max(q(3; s); q(3; d)) \\ q(3; s) &= 1 + \max(q(2; s); q(2; d)) \\ q(3; d) &= 0 + \max(q(1; s); q(1; d)); \end{aligned}$$

which implies

$$\begin{aligned} q(1; s) &= q(1; d) \\ q(1; d) &= 2 + \max(q(2; s); q(2; d)) \\ q(2; s) &= q(2; d) \\ q(2; d) &= \max(q(3; s); q(3; d)) \\ q(3; s) &= 1 + \max(q(2; s); q(2; d)) \\ q(3; d) &= \max(q(1; s); q(1; d)) \end{aligned}$$

Consider two solutions $q_1, q_2 \in Q_1$ defined as follows:

$$\begin{aligned} q_1(1; s) &= \frac{1}{2}; q_1(1; d) = \frac{3}{2}; \\ q_1(2; s) &= \frac{1}{2}; q_1(2; d) = \frac{1}{2}; \\ q_1(3; s) &= \frac{1}{2}; q_1(3; d) = \frac{1}{2}; \\ q_2(1; s) &= \frac{2}{3}; q_2(1; d) = \frac{2}{3}; \\ q_2(2; s) &= \frac{4}{3}; q_2(2; d) = \frac{1}{3}; \\ q_2(3; s) &= \frac{1}{3}; q_2(3; d) = \frac{2}{3}; \end{aligned}$$

The midpoint of q_1 and q_2 , $q = 0.5q_1 + 0.5q_2$, satisfies

$$\begin{aligned} q(1; s) &= \frac{1}{12}; q(1; d) = \frac{13}{12}; \\ q(2; s) &= \frac{11}{12}; q(2; d) = \frac{5}{12}; \\ q(3; s) &= \frac{1}{12}; q(3; d) = \frac{1}{12}; \end{aligned}$$

Note that

$$\frac{5}{12} = q(2; d) \notin \max(q(3; s); q(3; d)) = \frac{1}{12};$$

Therefore q does not satisfy the action-value optimality equation ((2)) and $q \notin Q_1$. Thus Q_1 is not convex. ■

B.5. Proof of Lemma 1

Part 1: (23) \Rightarrow (28):

Choose any solution of (23), v , and choose any $\hat{z} \in \mathcal{Z}$. Using Theorem 3.1 (e) by S&F,

$$(I - P)v \leq r - I\hat{\pi} + P^T v - v$$

(in weakly communicating SMDPs, the “ L ” set in Theorem 3.1 (e) is \emptyset).

Using the above inequality and Lemma 2.1 by S&F, we have

$$(I - P^T)v \leq Z(r - I\hat{\pi})$$

$$v \leq Z(r - I\hat{\pi}) + P^T v$$

Because v can be any element in \mathcal{V} ,

$$v(s) \leq \sup_{\hat{z} \in \mathcal{Z}} [Z(r - I\hat{\pi}) + P^T v](s); \forall s \in \mathcal{S} \quad (33)$$

By Theorem 4.1 (c) in S&F, there exists a deterministic optimal policy $\hat{\pi}$, which is apparently an element of $\hat{\mathcal{Z}}$ such that

$$v = Z(r - I\hat{\pi}) + P^T v$$

This result, along with (33) shows that

$$v(s) = \max_{\hat{z} \in \mathcal{Z}} [Z(r - I\hat{\pi}) + P^T v](s); \forall s \in \mathcal{S}$$

Part 1 is proven.

Part 2: (28) \Rightarrow (23):

Conversely, if v satisfies (28), define

$$v(s) \stackrel{\cdot}{=} \max_{s^0; r; I} \sum_{o} \rho(s^0; r; I | s; o) (r - I\hat{\pi} + v(s^0)); \forall s \in \mathcal{S} \quad (34)$$

We first show that $v \leq v$. For any $\pi \in \Pi$,

$$\begin{aligned}
 v(s) &= \max_{\sigma} \sum_{s^0; r; l} p(s^0; r; l | s; \sigma) (r + l \hat{r} + v(s^0)) \\
 &= [r + l \hat{r} + P v](s) && \text{because of the "max"} \\
 &= [r + l \hat{r} + P (Z (r + l \hat{r}) + P^T v)](s) && \text{by (28)} \\
 &= [(I + P Z) (r + l \hat{r}) + P^T v](s) && \text{rearranging terms} \\
 &= [(Z + P^T) (r + l \hat{r}) + P^T v](s) && \text{by Equation 2.2 in S\&F} \\
 &= [Z (r + l \hat{r}) + P^T v](s): && \text{by Theorem 3.1 (a) in S\&F} \\
 &= v(s) && \text{by (28)}
 \end{aligned}$$

We now show that $v \geq v$, which, together with $v \leq v$, implies $v = v$.

Let π be a deterministic policy achieving all J^{π} maxima in (34). Then we have

$$v - v = r + l \hat{r} + P v \tag{35}$$

$$(I - P)v - r - l \hat{r} : \tag{36}$$

Multiplying both sides by P^T , we have

$$P^T (I - P)v - P^T (r + l \hat{r}) :$$

The l.h.s. is 0 because $P^T P = P^T$. The r.h.s. is 0 because

$$P^T (r + l \hat{r}) - P^T (r + l r(\cdot)) = 0;$$

where we use $\hat{r} = r(\cdot)$; $\delta \in \Pi$. Combining the above two inequalities, we have

$$P^T (r + l \hat{r}) = 0: \tag{37}$$

Now we need to use the following lemma, which is essentially the same as Lemma 2.1 by S&F, except that the signs of inequalities are reversed. The proof follows the same arguments as those used in the proof of Lemma 2.1 by S&F.

Lemma 3 Fix any policy $\pi \in \Pi$, suppose that $P^T b = 0$ and $(I - P)x - b \geq 0$, then $(I - P^T)x - Z b \geq 0$.

Let $x = v$, and $b = r + l \hat{r}$, we see that $P^T b = 0$ because of (37) and $(I - P)x - b \geq 0$ because of (36). Using the above lemma, we have

$$v - Z (r + l \hat{r}) + P^T v :$$

Inserting this inequality to (35), we have

$$\begin{aligned}
 v &= [r \quad I \hat{r} + P v] \\
 &= (r \quad I \hat{r} + P [Z (r \quad I \hat{r}) + P^T v]) \\
 &= [(I + P Z)(r \quad I \hat{r}) + P^T v] \\
 &= [(Z + P^T)(r \quad I \hat{r}) + P^T v] && \text{by Equation 2.2 in S\&F} \\
 &= [Z (r \quad I \hat{r}) + P^T v] && \text{by Theorem 3.1 (a) in S\&F} \\
 &= \max_{\mathcal{Z}} [Z (r \quad I \hat{r}) + P^T v] \\
 &= v && \text{because } v \text{ satisfies (28):}
 \end{aligned}$$

Combining $v \leq v$ and $v \geq v$, we have $v = v$ and therefore $v(s) = v(s) = \max_{\sigma} \sum_{s^0:r;l} p(s^0; r; l | j; s; \sigma)(r \quad I \hat{r} + v(s^0))$; $\forall s \in S$. And thus v is a solution of (23).

B.6. Proof of Lemma 2

Choose $\epsilon > 0$. Because $v; v + \chi$ are both solutions of (23), using part (e)(2) of Theorem 3.1 by S&F, we have, $\forall s \in R$,

$$\begin{aligned}
 v(s) &= [r \quad I \hat{r} + P v](s) \\
 [v + \chi](s) &= [r \quad I \hat{r} + P (v + \chi)](s)
 \end{aligned}$$

Subtracting, we have $\chi(s) = [P \chi](s)$. Iteratively applying this equation, we have

$$\chi(s) = [P^d \chi](s) = \sum_{x_i} d_i \chi(x_i); \forall s \in R$$

where d is the unique stationary distribution of the d -th recurrent class of P . This proves (29).

Because $v + \chi$ is a solution of (23), by Lemma 1, $\forall s \in S$,

$$\begin{aligned}
 [v + \chi](s) &= \max_{\mathcal{Z}} [Z (r \quad I \hat{r}) + P^T [v + \chi]](s) \\
 \chi(s) &= \max_{\mathcal{Z}} [Z (r \quad I \hat{r} + P^T v - v) + P^T \chi](s) && \text{by Equation 2.2 in S\&F} \\
 &= \max_{\mathcal{Z}} [Z b_{v_i} + P^T \chi](s) \\
 &= \max_{\mathcal{Z}} [Z b_{v_i}](s) + \sum_{s^0 \in R} P^T(s; s^0) \chi(s^0)
 \end{aligned}$$

Using the above equality and (29), (30) and (31) hold.

B.7. Uniqueness of the Zero Solution

The other result we will need to use to show the convergence of the four algorithms introduced in the previous section is the following one. With this result, the stability of the algorithms can be established using the result by Borkar and Meyn (2000) (see also, Section 3.2 by Borkar 2009).

Lemma 4 *If an SMDP is weakly communicating and all rewards are 0, 0 is the only element in Q_1 .*

Proof Given a weakly communicating SMDP, by Lemma 1, any solution of the state-value optimality equation ((23)) satisfies

$$v(s) = \sum_{s^j \in R} P^1(s; s^j) v(s^j); \forall s \in R; \pi \in \hat{\Pi};$$

where R is defined right after Lemma 1. Also, because all rewards are 0, R is the closed communicating class and $\hat{\Pi} = \hat{\Pi}^R$ contains all stationary policies.

Pick an arbitrary $v \in \mathbb{R}^S$ and an arbitrary policy $\pi \in \hat{\Pi}$. For each recurrent class C under P , we have, by Lemma 1, $\forall s \in C, v(s) = \sum_{s^j \in C} d^C(s^j) v(s^j)$, where d^C denotes the stationary distribution of π in the recurrent class C . The r.h.s. only involves class C because starting from a state $s \in C$ the MDP can not leave C . Because the choice of s is arbitrary, $\min_{s \in C} v(s) = \sum_{s^j \in C} d^C(s^j) v(s^j) = \min_{s \in C} v(s)$. Thus $\sum_{s^j \in C} d^C(s^j) v(s^j) = \min_{s \in C} v(s)$. In addition, because $d^C(s) > 0$ for all $s \in C$, $v(s^j) = v(s); \forall s; s^j \in C$.

Now for any $s; s^j \in R$, there must exist a $\pi \in \hat{\Pi} = \hat{\Pi}^R$ such that there is a path from s to s^j and a path from s^j to s , because $s; s^j$ are in the same communicating class. Therefore $s; s^j$ are in the same recurrent class under P . Thus we conclude that $v(s) = v(s^j)$. Therefore $\forall s; s^j \in R, v(s) = v(s^j)$. The transient states values are uniquely determined by values of states in R . And in this case they are all equal to the values of states in the communicating class because all rewards are zero. Thus the solution set of v in the state-value optimality equation ((23)) is $fce: \forall c \in \mathbb{R}$.

Now consider the solution set of the option-value optimality equation ((11)). Let $v(s) = \max_o q(s; o)$, then (11) transforms to the state-value optimality equation. Therefore for any two solutions of q in (11), q_1 and q_2 , $\max_o q_1(s; o) = \max_o q_2(s; o) + ce$ for some c . Furthermore, let q be any solution of q , $q(s; o_1) = q(s; o_2); \forall s \in S; o_1; o_2 \in O$ because $\forall s \in S; o \in O$:

$$\begin{aligned} q(s; o) &= \sum_{s^j} p(s^j; r; j; s; o) \max_{o^j} q(s^j; o^j) \\ &= \sum_{s^j; r} p(s^j; r; j; s; o) \max_{o^j} q(s; o^j) \\ &= \max_{o^j} q(s; o^j); \end{aligned}$$

Thus the solution set of q in the option-value optimality equation ((11)) is also $fce: \forall c \in \mathbb{R}$. Given (22), $cu = cf(e) = f(ce) = \hat{r} = 0$ and $u > 0$ implies that $c = 0$. Therefore 0 is the unique solution of q . The lemma is proved. \blacksquare

B.8. General RVI Q

We now start to present the General RVI Q algorithm.

Let $I = \{1; 2; \dots; k\}$ where k is a positive integer. Consider solving $r \in \mathbb{R}$ and $q \in \mathbb{R}^{|I|}$ in following equation

$$r(i) = r + g(q)(i) \quad q(i) = 0; \forall i \in I \quad (38)$$

where $r \in \mathbb{R}^{j|j}$ is any fixed $j|j$ -dim vector, and $g : \mathbb{R}^{j|j} \rightarrow \mathbb{R}^{j|j}$ satisfies [Assumption 6](#).

We now consider an algorithm solving (38). This algorithm maintains an $j|j$ -dim vector of estimates $Q \in \mathbb{R}^{j|j}$, and updates Q using

$$Q_{n+1}(i) \stackrel{\dot{=}}{=} Q_n(i) + \frac{1}{(n;i)} (R_n(i) - F_n(Q_n) + G_n(Q_n)(i) - Q_n(i)) \mathbf{1}_{i \in Y_n}; \quad (39)$$

where

1. $\{Y_n\}$ is the “update schedule” – it is a set-valued process taking values in the set of nonempty subsets of I with the interpretation: $Y_n = \{i : i^{\text{th}} \text{ component of } Q \text{ was updated at time } n\}$,
2. $(n;i) \stackrel{\dot{=}}{=} \sum_{k=0}^n \mathbf{1}_{i \in Y_k}$, where $\mathbf{1}$ is the indicator function (i.e., $(n;i) =$ the number of times the i component was updated up to step n),
3. $\{c_n\}$ is a step-size sequence,
4. $\{R_n\}$ is a sequence of i.i.d. random vectors satisfying $E[R_n] = r; \forall n = 0; 1; \dots$,
5. for any $Q \in \mathbb{R}^{j|j}$, $\{G_n(Q)\}$ is a sequence of i.i.d. random vectors satisfying $E[G_n(Q)(i)] = g(Q)(i); \forall i \in I; n = 0; 1; \dots$,
6. for any $Q \in \mathbb{R}^{j|j}$, $\{F_n(Q)\}$ is a sequence of i.i.d. random variables satisfying $E[F_n(Q)] = f(Q); \forall n = 0; 1; \dots$ where $f : \mathbb{R}^{j|j} \rightarrow \mathbb{R}^{j|j}$ is a function satisfying [Assumption 4](#),
7. $\{c_n\}$ is a sequence of random vectors of size $j|j$.

We need the following assumptions in addition to Assumptions 1–4.

Assumption 6 1) g is a max-norm non-expansion, 2) g is a span-norm non-expansion, 3) $g(x + ce) = g(x) + ce$ for any $c \in \mathbb{R}; x \in \mathbb{R}^{j|j}$, 4) $g(cx) = cg(x)$ for any $c \in \mathbb{R}; x \in \mathbb{R}^{j|j}$.

Let

$$M_{n+1} \stackrel{\dot{=}}{=} R_n - r + G_n(Q_n) - g(Q_n) - (F_n(Q_n) - f(Q_n))e; \quad (40)$$

Let $\mathcal{F}_n \stackrel{\dot{=}}{=} (\mathcal{Q}_0; \mathcal{M}_1; \mathcal{M}_2; \dots; \mathcal{M}_n)$ be the increasing family of σ -fields. By the above construction, $\{M_{n+1}\}$ is a martingale difference sequence w.r.t. \mathcal{F}_n . That is $E[M_{n+1} | \mathcal{F}_n] = 0$ a.s., $n \geq 0$.

Assumption 7 For $n \geq 0; 1; 2; \dots$, $E[kR_n - rk^2 | \mathcal{F}_n] \leq K$, $E[kG_n(Q) - g(Q)k^2 | \mathcal{F}_n] \leq K(1 + kQk^2)$ for any $Q \in \mathbb{R}^{j|j}$, and $E[kF_n(Q) - f(Q)ek^2 | \mathcal{F}_n] \leq K(1 + kQk^2)$ for any $Q \in \mathbb{R}^{j|j}$ for a suitable constant $K > 0$.

We make the following assumption on $\{c_n\}$.

Assumption 8 $E[kc_nk^2 | \mathcal{F}_n] \leq K(1 + kc_nk^2)$ a.s.. Further, c_n converges to 0 a.s..

Assumption 9 Equation 38 has a unique solution of r . That is, there exists a pair $r \in \mathbb{R}; q \in \mathbb{R}^{I \times J}$ satisfying (38) and if both $r_1; q_1$ and $r_2; q_2$ are solutions (38), $r_1 = r_2$.

Denote the unique solution of r by $r_\#$.

Define $\mathcal{Q}_\#$ to be the set of $q \in \mathbb{R}^{I \times J}$ satisfying (38) and

$$r_\# = f(q): \quad (41)$$

Assumption 10 $\mathcal{Q}_\#$ is non-empty, bounded, closed, and connected.

Assumption 11 If $r(i) = 0; \forall i \in I$, then 0 is the only element in $\mathcal{Q}_\#$.

Assumption 5 is assumed in Section 2 of Wan et al. (2021b). Assumptions 1–8 are the same as Assumptions A.1–A.7 by Wan et al. (2021b). Assumptions 9–11 replace Assumption A.8 by Wan et al. (2021b).

Theorem 6 Under Assumptions 1–11, General RVI Q ((39)) converges, almost surely, Q_n to $\mathcal{Q}_\#$ and $f(Q_n)$ to $r_\#$.

Proof The proof would by a large degree repeats that of Theorem A.1 by Wan et al. (2021). For simplicity, we only highlight modifications.

Both our proof and the proof of Theorem A.1 study two ordinary differential equations (ODEs):

$$\dot{y}_t \doteq T_1(y_t) \quad y_t: \quad (42)$$

$$\dot{x}_t \doteq T_2(x_t) \quad x_t = T_1(x_t) \quad x_t + (r_\# \quad f(x_t)) e; \quad (43)$$

where

$$\begin{aligned} T_1(Q)(i) &\doteq r(i) + g(Q)(i) \quad r_\#; \\ T_2(Q)(i) &\doteq r(i) + g(Q)(i) \quad f(Q) \\ &= T_1(Q)(i) + (r_\# \quad f(Q)): \end{aligned}$$

Lemmas A.1 and A.3 by Wan et al. (2021) and their proofs still hold (by replacing r_1 with $r_\#$). Lemma A.2 by Wan et al. (2021) is replaced by the following one. The proof follows the same arguments as those for Lemma A.2.

Lemma 5 The set of equilibrium points of (43) is $\mathcal{Q}_\#$.

Lemma A.4 by Wan et al. (2021) is replaced with the following two lemmas.

Lemma 6 If $r(i) = 0; \forall i \in I$, 0 is the globally asymptotically stable equilibrium for (43).

Proof We have assumed that 0 is the only element in $\mathcal{Q}_\#$ in Assumption 11 and 0 is thus the unique equilibrium of (43). The rest follows the proof of Lemma A.4 by Wan et al. (2021) by replacing q_1 with 0. ■

Lemma 7 $Q_{\#}$ is a compact connected internally chain transitive invariant set for the ODE (43). Furthermore, any set that contains points not in $Q_{\#}$ is not a compact connected internally chain transitive invariant set for the ODE (43).

Proof According to Assumption 10, $Q_{\#}$ is closed and bounded and is thus compact. Assumption 10 also assumes that $Q_{\#}$ is connected. $Q_{\#}$ is an internally chain transitive invariant set for the ODE (43) because every element in $Q_{\#}$ is an equilibrium point of the ODE and $Q_{\#}$ is connected.

If a set contains a point $x \in \mathbb{R}^{|J|}$ that is not in $Q_{\#}$, this set can not be internally chain transitive because x is "transient" and the trajectory of (43) can not be arbitrarily close to x at some arbitrarily large time step. ■

The last step is to show convergence of a synchronous version of (39)

$$Q_{n+1}(i) \doteq Q_n(i) + \frac{1}{n} (R_n(i) - F_n(Q_n) + G_n(Q_n)(i) - Q_n(i)) \quad (44)$$

This result, together with Assumptions 2 and 3, guarantees convergence of the asynchronous update ((44)) by applying results from Section 7.4 by Borkar (2009)

Lemma 8 Equation 44 converges a.s. Q_n to $Q_{\#}$ as $n \rightarrow \infty$.

Proof The proof essentially follows that of Lemma A.5 with two changes. First, using Lemma 6 we can show that the ODE $\dot{x}_t = h_1(x_t) = g(x_t) - f(x_t)$ has the origin as the unique globally asymptotically stable equilibrium. Second, the proof of Lemma A.5 uses Borkar's (2009) Theorem 2, which proves that Q_n converges to a (possibly sample path dependent) compact connected internally chain transitive invariant set of $\dot{x}_t = h(x_t)$ where

$$h(Q_n)(i) \doteq r(i) - f(Q_n) + g(Q_n)(i) - Q_n(i)$$

Lemma 7 and Theorem 2 by Borkar (2009) together imply that Q_n must converge to $Q_{\#}$. ■

B.9. Verifying Assumption 9–Assumption 11

When casting General RVI Q to RVI Q-learning, (38) becomes (2). It is then clear that Assumption 9 satisfies because the associated MDP is communicating, Assumption 10 holds because of Theorem 5 and Assumption 11 holds because of Lemma 4.

When casting General RVI Q to Differential Q-learning, (38) becomes (2) for an MDP with the same transition dynamics as the original one and all rewards being shifted by a constant that depends on initial action-value estimate Q_0 and reward rate estimate R_0 . It is clear that this new MDP, just like the original one, is communicating. Therefore Theorem 5 and Lemma 4 hold and thus Assumption 9–Assumption 11 hold.

When casting General RVI Q to inter-option Differential Q-learning, (38) becomes (11) for an SMDP with the same transition dynamics as the original one and the reward of each state-option pair

being shifted in proportional to its expected duration. Again the resulting SMDP is communicating and therefore [Theorem 5](#) and [Lemma 4](#) hold and thus [Assumption 9–Assumption 11](#) hold.

When casting General RVI Q to intra-option Differential Q-learning, (38) becomes (12), which is the same as (11) by [Proposition 2](#) for the reward shifted SMDP introduced in the previous paragraph.

B.10. Convergence of Reward Rates of Greedy Policies

Lemma 9 *Assume that the SMDP is weakly communicating, suppose Q_n converges to Q_1 almost surely, let π_n be a greedy policy w.r.t. Q_n , $r(\pi_n; s) \rightarrow \hat{r}$ almost surely.*

Proof We will need the following lemma for the proof.

Lemma 10 *For any $m \geq 1; 2; 3; \dots; g$, $a; b \in \mathbb{R}^m$, $b > 0$, and for any $p \in \mathbb{R}^m$ such that $p > 0$, $\sum_s p(s) = 1$,*

$$\frac{p > a}{p > b} = \min_s \frac{a(s)}{b(s)};$$

$$\frac{p > a}{p > b} = \max_s \frac{a(s)}{b(s)};$$

Proof For any $s \in \{1; 2; 3; \dots; mg\}$

$$\frac{a(s)}{b(s)} = \min_{s'} \frac{a(s')}{b(s')}$$

$$a(s) = b(s) \min_{s'} \frac{a(s')}{b(s')}$$

$$p(s)a(s) = p(s)b(s) \min_{s'} \frac{a(s')}{b(s')}$$

Therefore

$$p > a = p > b \min_{s'} \frac{a(s')}{b(s')}$$

By our assumptions on b and p , $p > b > 0$, we have

$$\frac{p > a}{p > b} = \min_{s'} \frac{a(s')}{b(s')};$$

$\frac{p > a}{p > b} = \max_s \frac{a(s)}{b(s)}$ can be shown in the same way. ■

Given that Q_n converges to Q_1 , consider $r(\pi_n; s)$ where π_n is a greedy policy w.r.t. Q_n . We show that $r(\pi_n; s)$ converges to \hat{r} for all $s \in S$.

For any $\epsilon > 0$, let P^ϵ denote the $(s|j, j|o, s|j, j|o)$ transition probability matrix under policy π^ϵ . That is,

$$P^\epsilon(s; o; s^j; o^j) \doteq \prod_{r;l} p(s^j; r; l | j; s; o) \quad (o^j | j; s^j) \quad (45)$$

Let P^∞ be the *limiting matrix* of P^ϵ , which is the Cesaro limit of the sequence $\{P^\epsilon\}_{\epsilon=1}^\infty$:

$$P^\infty \doteq \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=0}^{n-1} P^i$$

Because S is finite, the Cesaro limit exists and P^∞ is a stochastic matrix (has row sums equal to 1). Let $r(s; o) \doteq \sum_{s^j; r; l} p(s^j; r; l | j; s; o) r$ be the one-stage option reward and $l(s; o) \doteq \sum_{s^j; r; l} p(s^j; r; l | j; s; o) l$ be the one-stage option length. Let $r(\pi; s; o) \doteq \sum_{s^j; r; l} p(s^j; r; l | j; s; o) r(\pi; s^j)$ be the reward rate of policy π starting from $s; o$. By part (a) of Theorem 11.4.1 in Puterman (1994),

$$r(\pi; s; o) = \frac{P_n^\infty r(s; o)}{P_n^\infty l(s; o)}$$

Thus we have,

$$\begin{aligned} r(\pi; s; o) &= \frac{P_n^\infty r(s; o)}{P_n^\infty l(s; o)} \\ &= \frac{P_n^\infty (r + P_n Q_n - Q_n)(s; o)}{P_n^\infty l(s; o)} \\ &= \min_{s^j; o^j} \frac{(r + P_n Q_n - Q_n)(s^j; o^j)}{l(s^j; o^j)} \\ &= \min_{s^j; o^j} \frac{(TQ_n - Q_n)(s^j; o^j)}{l(s^j; o^j)}; \end{aligned}$$

where the inequality holds because of [Lemma 10](#), and $TQ_n(s; o) \doteq r(s; o) + \sum_{s^j; r; l} p(s^j; r; l | j; s; o) Q_n(s^j; o^j)$.

Because the SMDP is weakly communicating, consider a deterministic optimal policy π^* , $r(\pi^*; s) = \hat{r}; \forall s \in S$. Therefore $r(\pi^*; s; o) = \hat{r}$. Now we have, for any $s; o$,

$$\begin{aligned} \hat{r} &= r(\pi^*; s; o) \\ &= \frac{P_n^\infty r(s; o)}{P_n^\infty l(s; o)} \\ &= \frac{P_n^\infty (r + P_n Q_n - Q_n)(s; o)}{P_n^\infty l(s; o)} \\ &= \max_{s^j; o^j} \frac{(r + P_n Q_n - Q_n)(s^j; o^j)}{l(s^j; o^j)} \\ &= \max_{s^j; o^j} \frac{(r + P_n Q_n - Q_n)(s^j; o^j)}{l(s^j; o^j)} \\ &= \max_{s^j; o^j} \frac{(TQ_n - Q_n)(s^j; o^j)}{l(s^j; o^j)}; \end{aligned}$$

The first inequality holds because of [Lemma 10](#) and the second inequality holds because q_n is a greedy policy w.r.t. Q_n .

With the above results, and that $\hat{r}(n; s; o) = r(n; s; o)$, we have, for any $s; o$,

$$\min_{s^l; o^l} \frac{(TQ_n - Q_n)(s^l; o^l)}{l(s^l; o^l)} = r(n; s; o) = \hat{r}(n; s; o) = \max_{s^l; o^l} \frac{(TQ_n - Q_n)(s^l; o^l)}{l(s^l; o^l)}$$

Therefore,

$$\max_{s; o} \hat{r}(n; s; o) = \text{sp} \left(\frac{TQ_n - Q_n}{l} \right);$$

where $\text{sp}(x) = \max_i x(i) - \min_i x(i)$ denotes the span of vector x .

Because $Q_n \rightarrow Q_1$ a.s., every point q in Q_1 satisfies $\text{sp}((Tq - q) = l) = 0$ because $(Tq - q) = l = \hat{r} e$ by (11). In addition, $\text{sp}((TQ_n - Q_n) = l)$ is a continuous function of Q_n , by continuous mapping theorem, $\text{sp}((TQ_n - Q_n) = l) \rightarrow 0$ a.s.. Therefore we conclude that $r(n; s; o) \rightarrow \hat{r}(s; o)$. By definition, $r(n; s) = \sum_{o \in \mathcal{O}(s)} p(o|s) r(n; s; o)$. Therefore $r(n; s) \rightarrow \hat{r}(s)$. ■

Appendix C. Empirical Results

In this section, we empirically verify our convergence results of Differential Q-learning and RVI Q-learning by showing the dynamics of estimated action values of the two algorithms in a communicating MDP and a weakly communicating MDP.

We first consider the communicating MDP shown at the bottom of [Figure 1](#). For this MDP, we apply Differential Q-learning with initial action values 0, initial reward rate estimate 3, and $\gamma = 1$. The behavior policy chooses action solid with probability 0.8 and action dashed with probability 0.2 for both two states. The stepsize is 0.1. We performed 10 runs for each algorithm. Each run starts from state 1 and lasts for 1000 steps. Every 10 steps, we recorded the estimated action values and plotted the higher action value for each state in the figure. [Figure 4](#)(left) shows the evolution of these action values. In the right panel of the same figure, we also show it using RVI Q-learning with action values being initialized with 0 and $f(q) = q(1; \text{dashed})$. A more detailed explanation is provided in the figure's caption. It can be seen that for both algorithms, 1) for each run, the estimated action-value function converged to a point in the solution set (the black line segments), and 2) for different runs, the estimated action values generally converged to different points in the solution set.

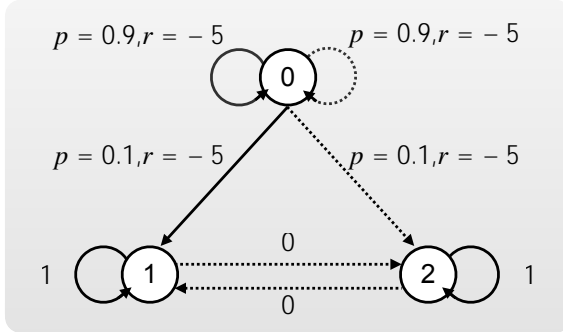


Figure 3: A weakly communicating MDP modified from the communicating MDP shown at the bottom of [Figure 1](#) by adding a transient state 0.

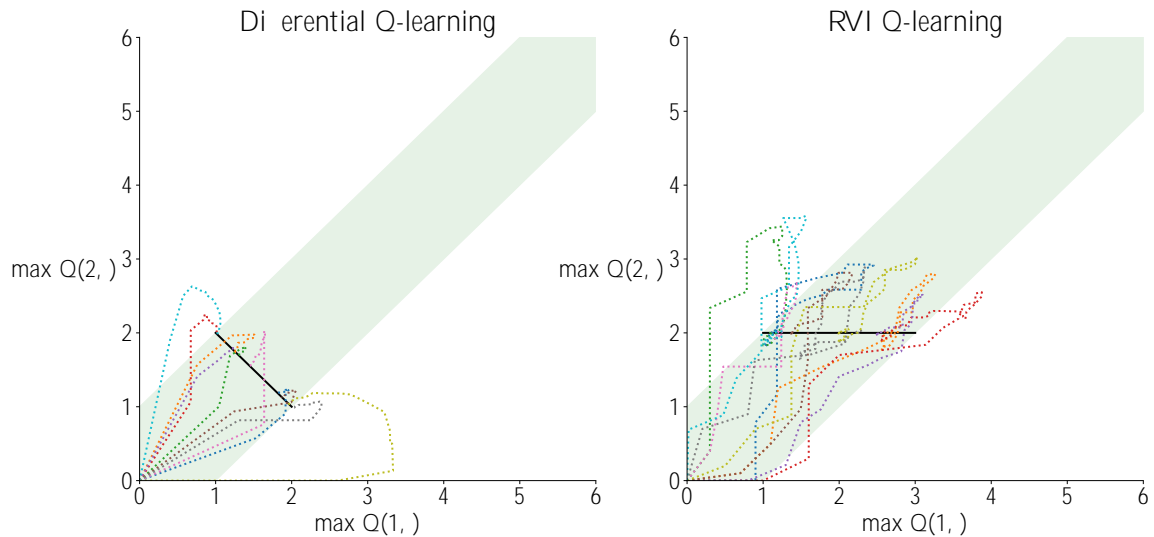


Figure 4: Evolution of estimated action values of Differential Q-learning and RVI Q-learning in the communicating MDP shown at the bottom of Figure 1. In each figure, the x and y axes are the higher estimated action value at state 1 and 2, respectively. The light green region marks the solution set of the action-value optimality equation (Equation (2)), and the black line segment marks the solution set Q_1 . Each colored dotted trajectory marks the evolution of the estimated action values. Each trajectory starts from zero point and ends at some point on the black line segment.

We also applied both of the two algorithms, with the same parameter settings and initialization in a weakly communicating MDP (Figure 3), which is just the communicating MDP plus a transient state. In the transient state, taking both solid and dashed actions stays at the transient state with probability 0.9. The MDP moves to state 1 with probability 0.1 given action solid and to state 2 with probability 0.1 given action dashed. The reward starting from state 0 is always 5. The starting state is 0. Because the agent could spend different amounts of time in the transient state for different runs, the agent may enter the communicating set, which contains states 1 and 2, with different action values associated with state 0.

The solution set of Differential Q-learning depends on the action values associated with the transient states when entering the communicating class. Therefore in the figure, the points that the estimated action-value function converged to, corresponding to different runs, are not in a line. Nevertheless, the estimated action-value function in all runs converged to the green region, which corresponds to the solution set of the action-value optimality equation.

On the other hand, The solution set of RVI Q-learning with the choice of the reference function $f(q) = q(1; \text{dashed})$ does not depend on the action values associated with states in the communicating class when entering the class. Therefore the solution set did not vary across different runs. Note that if we chose $f(q) = q(0; \text{dashed})$, then again the solution set of RVI Q-learning has that dependence.

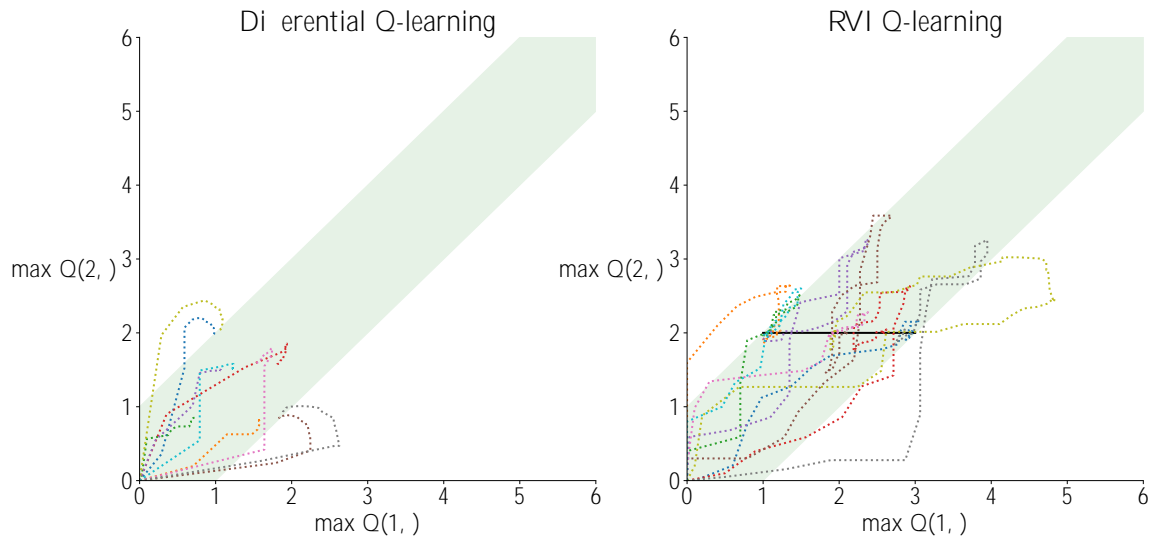


Figure 5: Evolution of estimated action values of Differential Q-learning and RVI Q-learning in the weakly communicating MDP shown in Figure 3.

Appendix D. Gosavi’s (2004) Convergence Result Is Incorrect

The convergence result of Gosavi’s proposed algorithm is presented in Theorem 2 of his paper. In the proof of the theorem, they used Borkar’s two-time scale stochastic approximation result to prove the convergence of the proposed algorithm. Specifically, they argued that their algorithm is a special case of the general class of algorithms considered in Borkar’s result. As Gosavi quotes, “Note that the Eqs. (48) and (49) for SMDPs form a special case of the general class of algorithms (29) and (30) analyzed using the lemma given in Section 5.1.1. ” However, a closer look at these equations shows that equation (49) is not a special case of equation (30). Note that because γ^k is a scalar, y^k only has one element and thus the f function in equation (30) does not vary across different state-option pairs. However, this is not true for the f function in equation (49). It appears to us that there is no simple fix for this issue.