

Multi-Objective Reinforcement Learning using Sets of Pareto Dominated Policies

Hasan Mert Gökalp

December 12, 2022

1 Introduction

2 Reinforcement Learning

- Background
- Single-objective
- Multi-Objective

3 Pareto Q-Learning

- Algorithm
- Set Evaluation Mechanisms
- Tracking a Policy

- Many real-life problems involve dealing with multiple objectives.
- No single optimal solution
- We want to obtain the best trade-off solutions(pareto front)
- Single-policy algorithms and multi-policy algorithms
- Reinforcement learning for multi-objective problems.

Markov Decision Process A Markov decision process is a 4-tuple (S, A, T_a, R_a) where:

- $S = (s_1, \dots, s_N)$ the state space
- $A = (a_1, \dots, a_r)$ the action set
- $T(s'|s, a)$ the transition probability
- $R(s, a)$ the expected immediate reward

Policies

- The goal is to learn a deterministic stationary policy π
- The state-dependent value function:

$$V^{\pi}(s) = E_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s\right]$$

- The action value function: $Q^{\pi}(s, a)$ stores the expected return starting from state s , taking action a , and thereafter following π again

Q-Learning

- The optimal Q^* -values:

$$Q^*(s, a) = R(s, a) + \gamma \sum_{s'} T(s' | s, a) \max_{a'} Q^*(s', a)$$

- Each entry contains a value for $\hat{Q}(s, a)$ which is the learner's current estimate about the actual value of $Q^*(s, a)$

$$\hat{Q}(s, a) \leftarrow (1 - \alpha_t) \hat{Q}(s, a) + \alpha_t \left(r + \gamma \max_{a'} \hat{Q}(s', a') \right)$$

Algorithm 1 Single-objective Q-learning algorithm

```
1: Initialize  $\hat{Q}(s, a)$  arbitrarily
2: for each episode  $t$  do
3:   Initialize  $s$ 
4:   repeat
5:     Choose  $a$  from  $s$  using a policy
6:     Take action  $a$  and observe  $s'$ 
7:      $\hat{Q}(s, a) \leftarrow (1 - \alpha_t) \hat{Q}(s, a) + \alpha_t \left( r + \gamma \max_{a'} \hat{Q}(s', a') \right)$ 
8:      $s \leftarrow s'$ 
9:   until  $s$  is terminal
10: end for
```

Multi Objective

- In the case of MORL, MDP provide a vector of rewards

$$R(s, a) = (R_1(s, a), \dots R_m(s, a))$$

- the state-dependent value function of a state s is vectorial

$$V^\pi(s) = E_\pi \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s \right]$$

- Since the environment now consists of multiple objectives, different policies can be optimal w.r.t. different objectives

Pareto front

Definition

A policy π_1 is said to strictly dominate another solution π_2 , that is $\pi_2 \prec \pi_1$, if each objective in V^{π_1} is not strictly less than the corresponding objective of V^{π_2} and at least one objective is strictly greater. A policy π is Pareto optimal if V^π either strictly dominates or is incomparable with the value functions of the other policies. The set of Pareto optimal policies is referred to as the Pareto front.

Single-policy Algorithms

Most approaches of reinforcement learning on multi-objective tasks rely on single-policy algorithms

- employ scalarization functions

$$v_w = f(v, w)$$

- scalar \hat{Q} -values are extended to \hat{Q} -vectors

$$\hat{Q}(s, a) = (\hat{Q}_1(s, a), \dots, \hat{Q}_m(s, a))$$

- scalarization of \hat{Q} vector for using traditional strategies:

$$\hat{S}Q(s, a) \leftarrow f(\hat{Q}(s, a), w)$$

Scalarization Algorithms

Selecting an action in a certain state of the environment

Algorithm 2 Scalarized ϵ -greedy strategy

```
1: SQList  $\leftarrow \{\}$ 
2: for each action  $a$  in  $A$  do
3:    $v \leftarrow \hat{Q}_1(s, a), \dots, \hat{Q}_m(s, a)$ 
4:    $\hat{SQ}(s, a) \leftarrow f(v, w)$ 
5:   Append  $\hat{SQ}(s, a)$  to SQList
6: end for
7: return  $\epsilon$ -greedy(SQList)
```

Algorithm 3 Scalarized multi-objective Q-learning algorithm

```
1: Initialize  $\hat{Q}_o(s, a)$  arbitrarily
2: for each episode  $t$  do
3:   Initialize state  $s$ 
4:   repeat
5:     Choose action  $a$  from  $s$  using the policy derived from  $\hat{S}Q$ -values
6:     Take action  $a$  and observe  $s'$ 
7:      $a' \leftarrow \text{greedy}(s')$ 
8:     for each objective  $o$  do
9:        $\hat{Q}_o(s, a) \leftarrow (1 - \alpha_t) \hat{Q}_o(s, a) + \alpha_t (r + \gamma \hat{Q}_o(s', a'))$ 
10:    end for
11:     $s \leftarrow s'$ 
12:  until  $s$  is terminal
13: end for
```

- In contrast to single-policy MORL, multi-policy algorithms do not reduce the dimensionality of the objective space but aim to learn a set of optimal solutions at once
- dynamic programming (DP) function:

$$\hat{Q}_{set}(s, a) = \mathbf{R}(s, a) \oplus \gamma \sum_{s' \in \mathcal{S}} T(s' | s, a) V^{ND}(s')$$

- The ND operator is a function that removes all Pareto dominated elements of the input set and returns the set of non-dominated elements

$$V^{ND}(s') = ND\left(\bigcup_{a'} \hat{Q}_{set}(s', a')\right)$$

- Assumption: Problem is episodic
- It does not assume any model
- 3 set evaluation mechanisms
- The set-based bootstrapping problem (lack of correspondence)
- Learning Immediate And Future Reward Separately
- $\bar{\mathcal{R}}(s, a) :=$ average observed immediate reward vector of (s, a)
 $ND_t(s, a) :=$ set of non-dominated vectors in the next state of s that is reached through action a at time step t .

$$\hat{Q}_{set}(s, a) \leftarrow \bar{\mathcal{R}}(s, a) \oplus \gamma ND_t(s, a)$$

Algorithm 4 Pareto Q-learning algorithm

- 1: Initialize $\hat{Q}_{set}(s, a)$'s as empty sets
 - 2: **for** each episode t **do**
 - 3: Initialize state s
 - 4: **repeat**
 - 5: Choose action a using a policy derived from the $\hat{Q}_{set}(s, a)$'s
 - 6: Take action a and observe state s' and reward vector r
 - 7: $ND_t(s, a) \leftarrow ND \left(\cup_{a'} \hat{Q}_{set}(s', a') \right)$
 - 8: $\bar{\mathcal{R}}(s, a) \leftarrow \bar{\mathcal{R}}(s, a) + \frac{r - \bar{\mathcal{R}}(s, a)}{n(s, a)}$
 - 9: $s \leftarrow s'$
 - 10: **until** s is terminal
 - 11: **end for**
-

Set Evaluation Mechanisms

- The hypervolume measure is a quality indicator that evaluates a particular set of vectorial solutions by calculating the volume with respect to its elements and a reference point
- The only quality indicator to be strictly monotonic with the Pareto dominance relation

Algorithm 5 Hypervolume Q_{set} evaluation

```
1: Retrieve current state  $s$ 
2:  $evaluations = \{ \}$ 
3: for each action  $a$  do
4:    $hv_a \leftarrow HV(\hat{Q}_{set}(s, a))$ 
5:   Append  $hv_a$  to evaluations
6: end for
7: return evaluations
```

Set Evaluation Mechanisms

- consider the number of Pareto dominating \hat{Q} -vectors of the \hat{Q} -set of each action
- the action that relates to the largest number of Pareto dominating \hat{Q} -vectors over all actions in s is selected.

Algorithm 6 Cardinality Q_{set} evaluation

```
1: Retrieve current state  $s$ 
2:  $allQs = \{ \}$ 
3: for each action  $a$  in  $s$  do
4:   for each  $\hat{Q}$  in  $\hat{Q}_{set}(s, a)$  do
5:     Append  $[a, \hat{Q}]$  to  $allQs$ 
6:   end for
7: end for  $NDQs \leftarrow ND(allQs)$ 
8: return  $NDQs$ 
```

Tracking a Policy

- one needs to select actions consistently in order to retrieve a desired policy

Algorithm 7 Track policy given the expected reward vector

```
1: target  $\leftarrow V^\pi(s)$ 
2: repeat
3:   for each  $a$  in  $A$  do
4:     Retrieve  $\bar{\mathcal{R}}(s, a)$ 
5:     Retrieve  $ND_t(s, a)$ 
6:     for each  $Q$  in  $ND_t(s, a)$  do
7:       if  $\gamma Q + \bar{\mathcal{R}}(s, a) = \text{target}$  then
8:          $s \leftarrow s' : T(s'|s, a) = 1$ 
9:          $\text{target} \leftarrow Q$ 
10:      end if
11:    end for
12:  end for
13: until  $s$  is not terminal
```