

ADVANCED ALGORITHMS

Temporal-difference (TD) learning control methods can be grouped into on-policy and off-policy methods. (Control methods also solve for prediction problems) A third class of these TD control methods are actor-critic.

1.1 Overview

The problem with the value-function approach to finding policies are those policies are deterministic, whereas optimal policies are often stochastic. The second drawback is a small change in the value of an action can change which action is selected for the given state. [3]

Actor-critic refers to a family of algorithms. First introduced by Barto et al. [1]

Example of multiple continuous action: [2]

1.2 Policy Gradient Methods

A policy is the probability that action a is taken at time t when the agent is in state s at time t with parameter θ . In other words,

$$\pi(a|s, \theta) = \Pr\{A_t = a | S_t = s, \theta_t = \theta\}$$

If the action space is discrete and not too large, an exponential soft-max distribution can be used in action selection,

$$\pi(a|s, \theta) \doteq \frac{e^{h(s,a,\theta)}}{\sum_b e^{h(s,b,\theta)}},$$

where

$$h(s, a, \boldsymbol{\theta}) = \boldsymbol{\theta}^\top \mathbf{x}(s, a),$$

where $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ and $\mathbf{x}(s, a) \in \mathbb{R}^{d'}$ is the feature vector.

Algorithm 1: Actor-critic for episodic problems

Input: a differentiable policy parameterization $\pi(a|s, \boldsymbol{\theta})$

Input: a differentiable state-value function parameterization $\hat{v}(s, \mathbf{w})$

Parameter(s): step sizes $\alpha^\theta > 0, \alpha^w > 0$

1 Initialise policy parameter $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ and state-value weights $\mathbf{w} \in \mathbb{R}^d$ (e.g. to 0) ;

2 **for** each episode **do**

3 Initialise S (first state of episode);

4 $I \leftarrow 1$;

5 **for** each step of episode or until S is terminal **do**

6 $A \sim \pi(\cdot|S, \boldsymbol{\theta})$;

7 Take action A , observe S', R ;

8 $\delta \leftarrow R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})$ (if S' is terminal, $\hat{v}(S', \mathbf{w}) \doteq 0$) ;

9 $\mathbf{w} \leftarrow \mathbf{w} + \alpha^w \delta \nabla \hat{v}(S, \mathbf{w})$;

10 $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha^\theta I \delta \nabla \ln \pi(A|S, \boldsymbol{\theta})$;

11 $I \leftarrow \gamma I$;

12 $S \leftarrow S'$;

13 **end**

14 **end**

Note that in actor-critic with traces with linear function approximation $\nabla \hat{v}(S, \mathbf{w})$ reduces to the feature vector $\mathbf{x}(S, a)$

Algorithm 2: Actor-critic with eligibility traces for episodic problems

Input: a differentiable policy parameterization $\pi(a|s, \theta)$ **Input:** a differentiable state-value function parameterization $\hat{v}(s, w)$ **Parameter(s):** step sizes $\alpha^\theta > 0, \alpha^w > 0$ **Parameter(s):** trace decay rates $\lambda^\theta \in [0, 1], \lambda^w \in [0, 1]$

```

1 Initialise policy parameter  $\theta \in \mathbb{R}^{d'}$  and state-value weights  $w \in \mathbb{R}^d$  (e.g. to 0) ;
2 for each episode do
3   Initialise  $S$  (first state of episode);
4    $z^\theta \leftarrow \mathbf{0}$  ( $d'$ -component eligibility trace vector) ;
5    $z^w \leftarrow \mathbf{0}$  ( $d$ -component eligibility trace vector) ;
6    $I \leftarrow 1$  ;
7   for each step of episode or until  $S$  is terminal do
8      $A \sim \pi(\cdot|S, \theta)$  ;
9     Take action  $A$ , observe  $S', R$ ;
10     $\delta \leftarrow R + \gamma \hat{v}(S', w) - \hat{v}(S, w)$  (if  $S'$  is terminal,  $\hat{v}(S', w) \doteq 0$ ) ;
11     $z^w \leftarrow \gamma \lambda^w z^w + \nabla \hat{v}(S, w)$  ;
12     $z^\theta \leftarrow \gamma \lambda^\theta z^\theta + I \nabla \ln \pi(A|S, \theta)$  ;
13     $w \leftarrow w + \alpha^w \delta z^w$  ;
14     $\theta \leftarrow \theta + \alpha^\theta \delta z^\theta$  ;
15     $I \leftarrow \gamma I$  ;
16     $S \leftarrow S'$  ;
17  end
18 end

```
