

# **MLMI7: Reinforcement Learning and Decision Making**

## **Function Approximation in RL**

Presented by David Krueger and Carl Rasmussen

Slides prepared by Milica Gašić

Lent Term

# Value-function approximation

- ▶ Representing value function as a table is not possible for large state spaces or continuous state spaces
- ▶ In this case value function can be a parameterised function with weight vector  $\theta \in \mathbb{R}^n$ :

$$V_{\pi}(s) \approx \hat{V}_{\pi}(s, \theta)$$

- ▶ The number of components of  $\theta$  is much less than the number of states ( $n \ll |\mathcal{S}|$ ), and changing one weight changes the estimated value of many states.
- ▶ When a single state is updated, the change generalises from that state to affect the values of many other states.

# Back-ups as input-output pairs

Value function estimation can be described as a series of back-ups:

Monte Carlo back-up  $s_t \mapsto R_t$

TD back-up  $s_t \mapsto r_t + \gamma \hat{V}(s_{t+1}, \theta)$

DP back-up  $s \mapsto E_\pi[r_t + \gamma \hat{V}(s_{t+1}, \theta) \mid s_t = s]$

Each can be seen as an example of the desired input-output behaviour of the value function.

This means that we can apply function approximation **but** only such that allows data to be obtained sequentially.

## Prediction Objective

In the case of approximation it is not possible to get the prediction in all states correct. Therefore, we produce a distribution over states which specifies how much we care about the error in each particular state  $d(s)$ .

The objective function is then **Mean Squared Value Error**:

$$MSVE(\theta) = \sum_s d(s) \left( V_{\pi}(s) - \hat{V}(s, \theta) \right)^2$$

Typically one chooses  $d(s)$  to be the fraction of time spent in  $s$  under the target policy  $\pi$  - *occupancy frequency*.

# On-policy distribution

- ▶  $h(s)$  denotes the probability that an episode begins in state  $s$
- ▶  $e(s)$  denotes the average time steps spent in state  $s$  in a single episode.

$$e(s) = h(s) + \sum_{\hat{s}} e(\hat{s}) \sum_{\hat{a}} \pi(\hat{a} \mid \hat{s}) p(s \mid \hat{s}, \hat{a})$$

This system of equations can be solved for the expected number of visits  $e(s)$  yielding the distribution:

$$d(s) = \frac{e(s)}{\sum_{s'} e(s')}$$

# Stochastic gradient descent

- ▶  $\hat{V}(s, \theta)$  is differentiable wrt  $\theta = (\theta_1, \theta_2, \dots, \theta_n)^\top$ .
- ▶  $\theta_t$  is updated at each of a series of discrete time steps,  $t = 0, 1, 2, 3, \dots$ .
- ▶ A sample  $s_t \mapsto V_\pi(s_t)$  consists of a (possibly random) state  $s_t$  and its true value under the policy  $\pi$ . We assume that states appear in examples with the same distribution,  $d(s)$ , over which we are trying to minimize the *MSVE*:

$$\begin{aligned}\theta_{t+1} &= \theta_t - \frac{1}{2}\alpha \nabla \left( V_\pi(s_t) - \hat{V}(s_t, \theta_t) \right)^2 \\ &= \theta_t + \alpha \left( V_\pi(s_t) - \hat{V}(s_t, \theta_t) \right) \nabla \hat{V}(s_t, \theta_t),\end{aligned}$$

$\alpha > 0$  is a step-size parameter.

- ▶ It's called “stochastic” because the update is done on only a single example, which has been selected stochastically.

# Target output

- ▶ In practise true value  $V_\pi(s_t)$  is not available during learning.
- ▶ Instead, we have  $s_t \mapsto U_t$  where  $U_t$  is a noisy estimate of  $V_\pi(s_t)$ . The general SGD method for state-value prediction is:

$$\theta_{t+1} = \theta_t + \alpha \left( U_t - \hat{V}(s_t, \theta_t) \right) \nabla \hat{V}(s_t, \theta_t)$$

- ▶ If  $U_t$  is an unbiased estimate ( $E[U_t] = V_\pi(s_t)$ ) for each  $t$ , then  $\theta_t$  is guaranteed to converge to a local optimum for decreasing  $\alpha$ .

# Prediction with function approximation

---

**Algorithm 1** Gradient Monte Carlo Algorithm for Approximating  $\hat{V} \approx V_\pi$

---

- 1: Input: the policy  $\pi$  to be evaluated
  - 2: Input: a differentiable function  $\hat{V}(s, \theta) : \mathcal{S} \times \mathbb{R}^n \rightarrow \mathbb{R}$
  - 3: Initialise  $\theta$
  - 4: **repeat**
  - 5:   Generate an episode  $s_0, a_0, r_1, \dots, r_T, s_T$  using  $\pi$
  - 6:   **for**  $t = 0, 1, \dots, T$  **do**
  - 7:      $\theta \leftarrow \theta + \alpha \left( R_t - \hat{V}(s_t, \theta) \right) \nabla \hat{V}(s_t, \theta)$
  - 8:   **end for**
  - 9: **until** convergence
-



## Semi-gradient methods

- ▶ If instead of MC, we are using TD or DP updates for prediction using SGD, ie we perform bootstrapping, they all depend on the current value of the weight vector  $\theta_t$
- ▶ This implies that they will be biased and that they will not produce a true gradient-descent method.
- ▶ They include only a part of the gradient and are called *semi-gradient methods*.

## Linear methods

One of the most important special cases of function approximation is that in which the approximate function,  $\hat{V}(s, \theta)$ , is a linear function of the weight vector,  $\theta$ :

$$\begin{aligned}\hat{V}(s, \theta) &= \theta^\top \cdot \phi(s) \\ &= \sum_i \theta_i \phi_i(s)\end{aligned}$$

$\phi = (\phi_1, \phi_2, \dots, \phi_n)^\top$ ,  $\phi_i(s) : \mathcal{S} \rightarrow \mathbb{R}$  are feature functions.

## Linear approximation

It is natural to use stochastic gradient descent updates with linear function approximation. The gradient of the approximate value function with respect to  $\theta$  in this case is:

$$\nabla \hat{V}(s, \theta) = \phi(s)$$

In linear case there is only one optimum.

## Semi-gradient TD update

$$\begin{aligned}\boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t + \alpha \left( r_{t+1} + \gamma \hat{V}(s_{t+1}, \boldsymbol{\theta}_t) - \hat{V}(s_t, \boldsymbol{\theta}_t) \right) \nabla \hat{V}(s, \boldsymbol{\theta}_t) \\ &= \boldsymbol{\theta}_t + \alpha \left( r_{t+1} + \gamma \boldsymbol{\theta}_t^\top \boldsymbol{\phi}_{t+1} - \boldsymbol{\theta}_t^\top \boldsymbol{\phi}_t \right) \boldsymbol{\phi}_t \\ &= \boldsymbol{\theta}_t + \alpha \left( r_{t+1} \boldsymbol{\phi}_t - \boldsymbol{\phi}_t (\boldsymbol{\phi}_t - \gamma \boldsymbol{\phi}_{t+1})^\top \boldsymbol{\theta}_t \right),\end{aligned}$$

where  $\boldsymbol{\phi}_t = \boldsymbol{\phi}(s_t)$ . Once the system has reached steady state, for any given  $\boldsymbol{\theta}_t$ , the expected next weight vector is:

$$E[\boldsymbol{\theta}_{t+1} \mid \boldsymbol{\theta}_t] = \boldsymbol{\theta}_t + \alpha(\mathbf{b} - A\boldsymbol{\theta}_t),$$

where

$$\begin{aligned}\mathbf{b} &= E[r_{t+1} \boldsymbol{\phi}_t] \\ A &= E[\boldsymbol{\phi}_t (\boldsymbol{\phi}_t - \gamma \boldsymbol{\phi}_{t+1})^\top]\end{aligned}$$

If the system converges to  $\boldsymbol{\theta}$ , then  $\mathbf{b} - A\boldsymbol{\theta} = 0$

## Least-Squares TD

TD with linear function approximation converges asymptotically, for appropriately decreasing step sizes, to the TD fixpoint:

$$\boldsymbol{\theta} = A^{-1}\mathbf{b}$$

$$A = E[\boldsymbol{\phi}_t(\boldsymbol{\phi}_t - \gamma\boldsymbol{\phi}_{t+1})^\top]$$

$$\mathbf{b} = E[R_{t+1}\boldsymbol{\phi}_t]$$

If this is so, then we don't need to compute the solution iteratively. Instead, we can calculate  $A$  and  $\mathbf{b}$  separately and then find the fixpoint.

# Least-Squares TD prediction

---

**Algorithm 2** Least-Squares TD

---

```
1: Input: policy  $\pi$ , features  $\phi(s) \in \mathbb{R}^n$ ,  $\phi(\text{terminal}) = 0$ 
2: Initialise  $\widehat{A}^{-1} = \epsilon^{-1} \mathbf{I}$ ,  $\hat{b} = \mathbf{0}$ 
3: repeat
4:   for each episode do
5:     Initialise  $s$  and obtain  $\phi$ 
6:     for each step do
7:       Choose  $a \sim \pi(\cdot | s)$ , take  $a$ , observe  $r, s'$ , obtain  $\phi'$ 
8:        $\mathbf{v} = \widehat{A}^{-1}(\phi - \gamma\phi')$ 
9:        $\widehat{A}^{-1} = \widehat{A}^{-1} - (\widehat{A}^{-1}\phi)\mathbf{v}^T / (1 + \mathbf{v}^T\phi)$ 
10:       $\hat{b} \leftarrow \hat{b} + r\phi$ 
11:       $\theta = A^{-1}\hat{b}$ 
12:       $s \leftarrow s', \phi \leftarrow \phi'$ 
13:   end for
14: end for
15: until convergence
```

---

# Properties of LSTD

- ▶ Complexity is  $O(n^2)$  vs  $O(n)$  for semi-gradient TD
- ▶ No step size parameter is required
- ▶  $\epsilon$ -greedy policy is used in the policy improvement step
- ▶ This requires setting  $\epsilon$ : if  $\epsilon$  is too small the sequence of inverses can vary wildly, and if  $\epsilon$  is too large then learning is slowed
- ▶ It never forgets which is problematic if the target policy changes as it does in reinforcement learning and generalised policy iteration.

# Summary

- ▶ Reinforcement learning systems must be capable of generalization if they are to be applicable to artificial intelligence or to large engineering applications.
- ▶ In parameterised function approximation the value function is parameterised by a weight vector  $\theta$
- ▶ To find a good weight vector we use a variation of stochastic gradient descent
- ▶ Good results can be obtained for semi-gradient methods in the special case of linear function approximation, in which the value estimates are weighted sum of features.
- ▶ LSTD is the most data-efficient linear TD prediction method, but has computational complexity  $O(n^2)$  for  $n$  features



# Policy gradient methods

- ▶ Policy gradient methods learn a parametrised policy that can select actions without needing to compute a value function
- ▶ Policy  $\pi$  is parametrised with  $\omega \in \mathbb{R}^n$

$$\pi(a \mid s, \omega) = p(a_t = a \mid s_t = s, \omega_t = \omega)$$

- ▶ Given a performance measure  $J(\omega)$  the gradient is

$$\omega_{t+1} = \omega_t + \alpha \nabla J(\omega_t)$$

- ▶  $J(\omega)$  is typically the value of the initial state  $V_{\pi(\omega)}(s_0)$ ,

# Policy approximation

- ▶ Stochastic policy
- ▶ Approximation method such that gradient  $\nabla_{\omega}\pi(a|s, \omega)$  exists and is finite
- ▶ We often use a *Gibbs policy*:

$$\pi(a|s, \omega) = \frac{\exp(\omega^T \psi(s, a))}{\sum_{a'} \exp(\omega^T \psi(s, a'))}$$

where  $\psi$  denotes parametrised the feature functions.

# Policy gradient theorem

$$\nabla J(\omega) = \sum_s d_\pi(s) \sum_a Q_\pi(s, a) \nabla_\omega \pi(a \mid s, \omega)$$

PROOF

# REINFORCE

- ▶ The policy gradient theorem gives us an exact expression for the gradient; all we need is some way of sampling whose expectation equals or approximates this expression.
- ▶ Notice that the right-hand side is a sum over states weighted by how often the states occurs under the target policy  $\pi$  weighted again by  $\gamma$  times how many steps it takes to get to those states.
- ▶ If we just follow  $\pi$  we will encounter states in these proportions, which we can then weight by  $\gamma$  to preserve the expected value.

# REINFORCE

$$\begin{aligned}\nabla J(\omega) &= E_{\pi} \left[ \gamma^t R_t \frac{\nabla_{\omega} \pi(a|s, \omega)}{\pi(a|s_t, \omega)} \right] \\ &= E_{\pi} \left[ \gamma^t R_t \nabla_{\omega} \log \pi(a|s_t, \omega) \right] \\ \omega_{t+1} &= \omega_t + \alpha \gamma^t R_t \nabla \log \pi(a|s_t, \omega)\end{aligned}$$

In case  $\pi$  is a Gibbs policy:

$$\nabla \log \pi(a|s, \omega) = \psi(s, a) - \sum_b \pi(b|s, \omega) \psi(s, b)$$

# REINFORCE, A Monte-Carlo Policy-Gradient Method

---

## Algorithm 3 REINFORCE

---

- 1: Input: a differentiable policy parameterization  $\pi(a|s, \omega)$ ,  $\alpha > 0$
  - 2: Initialise  $\omega$
  - 3: **repeat**
  - 4:   Generate an episode  $s_0, a_0, r_1, \dots, s_T, a_T$  following  $\pi(\cdot|\cdot, \omega)$
  - 5:   **for** each step  $t = 0, \dots, T$  **do**
  - 6:      $R_t \leftarrow$  return from step  $t$
  - 7:      $\omega \leftarrow \omega + \alpha \gamma^t R_t \nabla \log \pi(a|s_t, \omega)$
  - 8:   **end for**
  - 9: **until** convergence
-

# Summary

- ▶ Instead of parametrising value functions we can directly parametrise policy
- ▶ Policy gradient theorem states the value of the gradient
- ▶ An episodic Monte Carlo algorithm which estimates policy parameters using policy gradient theorem is REINFORCE algorithm

## Next lecture

- ▶ Actor-critic methods