

# TD-learning and $Q$ -learning

<http://bicmr.pku.edu.cn/~wenzw/bigdata2020.html>

Acknowledgement: this slides is based on Prof. Shipra Agrawal's lecture notes

# Outline

- 1 TD-learning
  - TD(0)-learning
  - TD( $\lambda$ )
- 2  $Q$ -learning (tabular)
  - The  $Q$ -learning method
  - Stochastic Approximation method
- 3  $Q$ -learning with function approximation
- 4 Deep Q-learning Networks (DQN)
- 5 Approximate dynamic programming
  - TD(0) and TD(1)
  - Fitted Value iteration

# TD-learning

- TD-learning is essentially approximate version of policy evaluation without knowing the model (using samples). Adding policy improvement gives an approximate version of policy iteration.
- Since the value of a state  $V^\pi(s)$  is defined as the expectation of the random return when the process is started from the given state  $s$ , an obvious way of estimating this value is to compute an average over multiple independent realizations started from the given state. This is an instance of the so-called **Monte-Carlo** method.

# TD-learning

- Unfortunately, the variance of the observed returns can be high, which means that the quality of the estimates will be poor. The Monte-Carlo technique is further difficult to apply if the system is not accessible through a simulator but rather estimation happens while actually interacting with the system.
- In this case, when the rewards obtained while actually taking the actions are used as feedback to learn in a closed loop, it might be impossible to reset the state of the system to some particular state.
- In this case, the Monte-Carlo technique cannot be applied without introducing some additional bias.
- Temporal difference (TD) learning (Sutton [1988]), which is one of the most significant ideas in reinforcement learning, is a method that can be used to address these issues.

# TD(0)-learning

- Policy evaluation is about estimating  $V^\pi(\cdot)$ , which by Bellman equations is equivalent to finding a stationary point of the following equations:

$$V^\pi(s) = \mathbb{E}_{a \sim \pi(s)} \left[ R(s, a, s') + \gamma \sum_{s'} P(s, a, s') V^\pi(s') \right], \forall s$$

- However, we need to estimate this using only observations  $r_t, s_{t+1}$  on playing some action at  $a_t$  current state  $s_t$ .
- Let the current estimate of value function for any state  $s$  is  $\hat{V}(s)$ . Let on taking action  $a_t = \pi(s_t)$  in the current state  $s_t, s_{t+1}$  is the observed (sample) next state. The predicted value function for the next state  $s_{t+1}$  is  $\hat{V}(s_{t+1})$ , giving another prediction of value function at state  $s_t$  (by one-lookahead using Bellman equations is)  $r_t + \gamma \hat{V}(s_{t+1})$ . Note that

$$\mathbb{E} [r_t + \hat{V}(s_{t+1}) | s_t, \hat{V}] = \mathbb{E}_{a \sim \pi(s_t)} \left[ R(s_t, a) + \sum_{s'} P(s_t, a, s') \hat{V}(s') | s_t, \hat{V} \right]$$

# TD(0)-learning

- From Bellman equations, we are looking for  $\hat{V}$  such that

$$\hat{V}(s_t) \approx r_t + \hat{V}(s_{t+1})$$

- Then, the **TD** method performs the following update to the value function estimate at  $s_t$ , moving it towards the new estimate:

$$\hat{V}(s_t) \leftarrow (1 - \alpha_t) \hat{V}(s_t) + \alpha_t (r_t + \gamma \hat{V}(s_{t+1}))$$

- Let  $\delta_t$  be the following gap:

$$\delta_t := r_t + \gamma \hat{V}(s_{t+1}) - \hat{V}(s_t)$$

- referred to as **temporal difference**, i.e., the difference between current estimate, and one-lookahead estimate. Then, the above also be written as:

$$\hat{V}(s_t) \leftarrow \hat{V}(s_t) + \alpha_t \delta_t \tag{1}$$

# SGD interpretation

- The above update can be interpreted as a stochastic gradient descent step for minimizing square loss (least squares method). Following supervised terminology, let  $\hat{V}(s_t)$  denotes a predictor. Let  $z_t = r_t + \gamma \hat{V}(s_{t+1})$  is an observation (also referred to as ‘target’). Then, the least squares method fits  $\hat{V}(s)$  for any  $s$  by minimizing squared loss:

$$\sum_{t:s_t=s} \frac{1}{2} (\hat{V}(s_t) - z_t)^2$$

- The loss function aims to minimize distance between prediction and observation.

# SGD interpretation

- To minimize the above loss function, one could use stochastic gradient descent. The gradient at  $t^{th}$  step with respect to  $\hat{V}(s_t)$  is  $(\hat{V}(s_t) - z) = -\delta_t$ . Thus, (1) is a gradient step with step size  $\alpha_t$ , which moves the prediction closer to the observations. This is also a popular rule in supervised learning called the LMS update rule (LMS stands for “least mean squares”), and is also known as the Widrow-Hoff learning rule.
- However, an important difference from use of squared error minimization supervised learning is that the ‘target’  $z_t = r_t + \gamma \hat{V}(s_{t+1})$  in the the above squared loss objective itself depends on the current estimate, thus the algorithm uses a form of ‘bootstrapping’.
- The gradient descent interpretation will be more significant in the large scale case when function approximation is used ( $\hat{V}$  will be estimated as a parametric function of state features).



# Tabular TD(0) method for policy evaluation

---

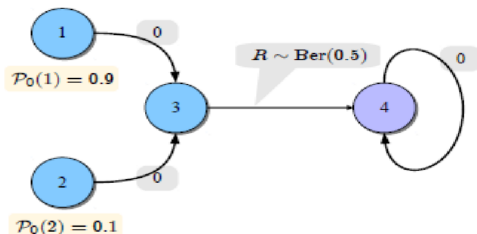
**Algorithm 1** Tabular TD(0) method for policy evaluation

---

- 1: Initialization: Given a starting state distribution  $D_0$ , policy  $\pi$ , the method evaluates  $V^\pi(s)$  for all states  $s$ .  
Initialize  $\hat{V}$  as an empty list/array for storing the value estimates.
  - 2: **repeat**
  - 3:   Set  $t = 1, s_1 \sim D_0$ . Choose step sizes  $\alpha_1, \alpha_2, \dots$ .
  - 4:   Perform TD(0) updates over an episode:
  - 5:   **repeat**
  - 6:     Take action at  $a_t \sim \pi(s_t)$ . Observe reward  $r_t$ , and new state  $s_{t+1}$ .
  - 7:      $\delta_t := r_t + \gamma \hat{V}(s_{t+1}) - \hat{V}(s_t)$
  - 8:     Update  $\hat{V}(s_t) \leftarrow \hat{V}(s_t) + \alpha_t \delta_t$
  - 9:      $t = t + 1$
  - 10:   **until** episode terminates
  - 11: **until** change in  $\hat{V}$  over consecutive episodes is small
-

# Monte Carlo method

Why use only 1-step lookahead to construct target  $z$ ? Why not lookahead entire trajectory (in problems where there is a terminal state, also referred to as episodic MDPs)?



[Szepesvari, 1999] In this example, all transitions are deterministic. The reward is zero, except when transitioning from state 3 to state 4, when it is given by a Bernoulli random variable with parameter 0.5. State 4 is a terminal state. When the process reaches the terminal state, it is reset to start at state 1 or 2. The probability of starting at state 1 is 0.9, while the probability of starting at state 2 is 0.1.

# Monte Carlo method

- The resulting method is referred to as Monte Carlo method, here for  $z$  a sample trajectory starting at  $s_t$  is used

$$z = \sum_{n=0}^{\infty} \gamma^n r_{t+n} =: \mathcal{R}_t$$

so that

$$\delta_t = z - \hat{V}(s_t) = \mathcal{R}_t - \hat{V}(s_t)$$

$$\hat{V}(s_t) = (1 - \alpha)\hat{V}(s_t) + \alpha\mathcal{R}_t$$

# TD(0) or Monte-Carlo?

*This example is taken from page 22 – 23, Szepesvari [1999]*

- First, let us consider an example when  $TD(0)$  converges faster. Consider the above undiscounted episodic MRP shown on the above figure.
- The initial states are either 1 or 2. With high probability the process starts at state 1, while the process starts at state 2 less frequently.
- Consider now how  $TD(0)$  will behave at state 2. By the time state 2 is visited the  $k^{th}$  time, on the average state 3 has already been visited  $10k$  times.
- Assume that  $\alpha_t = 1/(t + 1)$  (the TD updates with this step size reduce to averaging of target observations). At state 1 and 2, the target is  $\hat{V}(3)$  (since immediate reward is 0 and transition probability to state 3 is 1).

## TD(0) or Monte-Carlo?

- Therefore, whenever state 2 is visited the  $TD(0)$  sets its value as the average of estimates  $\hat{V}^t(3)$  over the time steps  $t$  when state 1 was visited (similarly for state 2). At state 3 the  $TD(0)$  update reduces to averaging the Bernoulli rewards incurred upon leaving state 3. At the  $k^{th}$  visit of state 2,  $\text{Var}(\hat{V}(3)) \simeq 1/(10k)$ . Clearly,  $\mathbb{E}[\hat{V}(3)] = 0.5$ . Thus, the target of the update of state 2 will be an estimate of the true value of state 2 with accuracy increasing with  $k$ .
- Now, consider the Monte-Carlo method. The Monte-Carlo method ignores the estimate of the value of state 3 and uses the Bernoulli rewards directly. In particular,  $\text{Var}(\mathcal{R}_t | s_t = 2) = 0.25$ , i.e., the variance of the target does not change with time.
- On this example, this makes the Monte-Carlo method slower to converge, showing that sometimes bootstrapping might indeed help.

## TD(0) or Monte-Carlo?

- To see an example when bootstrapping is not helpful, imagine that the problem is modified so that the reward associated with the transition from state 3 to state 4 is made deterministically equal to one.
- In this case, the Monte-Carlo method becomes faster since  $\mathcal{R}_t = 1$  is the true target value, while for the value of state 2 to get close to its true value,  $TD(0)$  has to wait until the estimate of the value at state 3 becomes close to its true value. This slows down the convergence of  $TD(0)$ .
- In fact, one can imagine a longer chain of states, where state  $i + 1$  follows state  $i$ , for  $i \in 1, \dots, N$  and the only time a nonzero reward is incurred is when transitioning from state  $N - 1$  to state  $N$ .
- In this example, the rate of convergence of the Monte-Carlo method is not impacted by the value of  $N$ , while  $TD(0)$  would get slower with  $N$  increasing.

# TD( $\lambda$ )

- TD( $\lambda$ ) is a "middle-ground" between  $TD(0)$  and Monte-Carlo evaluation.
- Here, the algorithm considers  $\ell$ -step predictions:

$$z_t^\ell = \sum_{n=0}^{\ell} \gamma^n r_{t+n} + \gamma^{\ell+1} \hat{V}(s_{t+\ell+1})$$

with temporal difference:

$$\begin{aligned}\delta_t^\ell &= z_t^\ell - \hat{V}(s_t) \\ &= \sum_{n=0}^{\ell} \gamma^n r_{t+n} + \gamma^{\ell+1} \hat{V}(s_{t+\ell+1}) - \hat{V}(s_t) \\ &= \sum_{n=0}^{\ell} \gamma^n (r_{t+n} + \gamma \hat{V}(s_{t+n+1}) - \hat{V}(s_{t+n})) \\ &= \sum_{n=0}^{\ell} \gamma^n \delta_{t+n}\end{aligned}$$

- In TD( $\lambda$ ) method, a mixture of  $\ell$ -step predictions is used, with weight  $(1 - \lambda)\lambda^\ell$  for  $\ell \geq 0$ . Therefore,  $\lambda = 0$  gives  $TD(0)$ , and  $\lambda \rightarrow 1$  gives Monte-Carlo method.  $\lambda > 1$  gives a multi-step method. To summarize, the  $TD(\lambda)$  update is given as:

$$\hat{V}(s_t) \leftarrow \hat{V}(s_t) + \alpha_t \sum_{\ell=0}^{\infty} (1 - \lambda) \lambda^\ell \delta_t^\ell = \hat{V}(s_t) + \alpha_t \sum_{n=0}^{\infty} \lambda^n \gamma^n \delta_{t+n}$$



# Policy improvement with TD-learning

TD-learning allows evaluating a policy. For using TD-learning for finding optimal policy, we need to be able to improve the policy. Recall policy iteration effectively requires evaluating  $Q$ -value of a policy, where  $Q^\pi(s, a) = R(s, a) + \gamma \sum_{s'} P^\pi(s, a, s') V^\pi(s')$ . With simple modification, TD-learning can be used to estimate  $Q$ -value of a policy. There, the updates would be replaced by:

$$\delta_t := r_t + \gamma \hat{Q}(s_{t+1}, \pi(s_{t+1})) - \hat{Q}(s_t, a_t)$$

$$\text{Update } \hat{Q}(s_t, a_t) \leftarrow \hat{Q}(s_t, a_t) + \alpha_t \delta_t$$

Then, the scheme for policy improvement is similar to policy iteration. Repeat the following until convergence to some policy:

- Use TD-learning to evaluate the policy  $\pi^k$ . The method outputs  $\hat{Q}^{\pi^k}(s, a), \forall s, a$
- Compute new 'improved policy'  $\pi^{k+1}$  as  $\pi^{k+1}(s) \leftarrow \arg \max_a \hat{Q}^{\pi^k}(s, a)$ .

# Outline

- 1 TD-learning
  - TD(0)-learning
  - TD( $\lambda$ )
- 2 *Q*-learning (tabular)
  - The *Q*-learning method
  - Stochastic Approximation method
- 3 *Q*-learning with function approximation
- 4 Deep Q-learning Networks (DQN)
- 5 Approximate dynamic programming
  - TD(0) and TD(1)
  - Fitted Value iteration

# $Q$ -learning (tabular)

- $Q$ -learning is a sample based version of  $Q$ -value iteration. This method attempts to directly find optimal  $Q$ -values, instead of computing  $Q$ -values of a given policy.
- Recall  $Q$ -value iteration: for all  $s, a$  update,

$$Q_{k+1}(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} P(s, a, s') \left( \max_{a'} Q_k(s', a') \right)$$

$Q$ -learning approximates these updates using sample observations, similar to TD-learning.

- In steps  $t = 1, 2, \dots$  of an episode, the algorithm observes reward  $r_t$  and next state  $s_{t+1} \sim P(\cdot, s_t, a_t)$  **for some action**  $a_t$ . It updates the  $Q$ -estimates for pair  $(s_t, a_t)$  as follows:

$$Q_{k+1}(s_t, a_t) = (1 - \alpha)Q_k(s_t, a_t) + \alpha \underbrace{\left( r_t + \gamma \max_{a'} Q_k(s_{t+1}, a') \right)}_{\text{target}}$$

# $Q$ -learning (tabular)

---

## Algorithm 2 Tabular $Q$ -learning method

---

- 1: Initialization: Given a starting state distribution  $D_0$ .  
Initialize  $\hat{Q}$  as an empty list/array for storing the  $Q$ -value estimates.
  - 2: **repeat**
  - 3:   Set  $t = 1$ ,  $s_1 \sim D_0$ . Choose step sizes  $\alpha_1, \alpha_2, \dots$ .
  - 4:   Perform  $Q$ -learning updates over an episode:
  - 5:   **repeat**
  - 6:     Take action at  $a_t$ . Observe reward  $r_t$ , and new state  $s_{t+1}$ .
  - 7:      $\delta_t := \left( r_t + \gamma \max_{a'} \hat{Q}(s_{t+1}, a') \right) - \hat{Q}(s_t, a_t)$
  - 8:     Update  $\hat{Q}(s_t, a_t) \leftarrow \hat{Q}(s_t, a_t) + \alpha_t \delta_t$
  - 9:      $t = t + 1$
  - 10:   **until** episode terminates
  - 11: **until** change in  $\hat{Q}$  over consecutive episodes is small
-

# How to select actions, the issue of exploration

- The convergence results (discussed below) for  $Q$ -learning will say that if all actions and states are infinitely sampled, learning rate is small, but does not decrease too quickly, then  $Q$ -learning converges. (Does not matter how you select actions, as long as they are infinitely sampled).
- One option is to select actions greedily according to the current estimate  $\max_a Q_k(s, a)$ . But, this will reinforce past errors, and may fail to sample and estimate  $Q$ -values for actions which have higher error levels. We may get stuck at a subset of (suboptimal) actions.
- Therefore, exploration is required. The  $\epsilon$ -greedy approach (i.e., with  $\epsilon$  probability pick an action uniformly at random instead of greedy choice) can ensure infinite sampling of every action, but can be very inefficient.

# How to select actions, the issue of exploration

- The same issue occurs in TD-learning based policy improvement methods. The choice of action is specified as the greedy policy according to the previous episode estimates. Without exploration this may not ensure that all actions and states are infinitely sampled.
- One option is to replace policy improvement step by greedy choice. That is, the policy improvement step will now compute the new 'improved policy'  $\pi^{k+1}$  as the randomized policy:

$$\pi^{k+1}(s) = \begin{cases} a_k^* := \arg \max_a \hat{Q}^{\pi^k}(s, a), & \text{with probability } 1 - \epsilon + \frac{\epsilon}{|A|} \\ a, & \text{with probability } \frac{\epsilon}{|A|}, a \neq a_k^* \end{cases}$$

- Then, in policy evaluation, this 'randomized policy' must be used.  
7 :  $\delta_t := r_t + \gamma \mathbb{E}_{a \sim \pi^{k+1}(s_{t+1})} [\hat{Q}(s_{t+1}, a)] - \hat{Q}(s_t, a_t)$   
8 : Update  $\hat{Q}(s_t, a_t) \leftarrow \hat{Q}(s_t, a_t) + \alpha_t \delta_t$

# Convergence theorem

## Theorem 1 (Watkins and Dayan [1992])

Given bounded rewards  $|r_t| \leq R$ , learning rates  $0 \leq \alpha_t < 1$ , and

$$\sum_{i=1}^{\infty} \alpha_{n^i(s,a)} = \infty, \sum_{i=1}^{\infty} (\alpha_{n^i(s,a)})^2 < \infty$$

then  $\hat{Q}^t(s, a) \rightarrow Q(s, a)$  as  $t \rightarrow \infty$  for all  $s, a$  with probability 1. Here,  $n^i(s, a)$  is the index of the  $i^{\text{th}}$  time the action  $a$  is tried in state  $s$ , and  $\hat{Q}^t(s, a)$  is the estimate  $\hat{Q}$  in round  $t$ .

- If all actions and states are infinitely sampled, learning rate is small, but does not decrease too quickly, then  $Q$ -learning converges. (Does not matter how you select actions, as long as they are infinitely sampled).
- The proof of this and many similar results in RL algorithms follow the analysis of a more general online learning/optimization method - the stochastic approximation method.

# Stochastic Approximation method

- The stochastic approximation (SA) algorithm essentially solves a system of (nonlinear) equations of the form

$$h(\theta) = 0$$

for unknown  $h(\cdot)$ , based on noisy measurements of  $h(\theta)$ .

- More specifically, consider a (continuous) function  $\mathbb{R}^d \rightarrow \mathbb{R}^d$ , with  $d \geq 1$ , which depends on a set of parameters  $\theta \in \mathbb{R}^d$ . Suppose that  $h(\theta)$  is unknown. However, for any  $\theta$  we can measure  $Z = h(\theta) + \omega$ , where  $\omega$  is some 0-mean noise. The classical SA algorithm (Robbins and Monro [1951]) is of the form

$$\begin{aligned}\theta_{n+1} &= \theta_n + \alpha_n Z_n \\ &= \theta_n + \alpha_n (h(\theta_n) + \omega_n), \quad n \geq 0\end{aligned}$$

- Since  $\omega_n$  is 0-mean noise, the stationary points of the above algorithm coincide with the solutions of  $h(\theta) = 0$ .



# Asynchronous version

- More relevant to the RL methods discussed here is the asynchronous version of the SA method. In the asynchronous version of SA method, we may observe only one coordinate (say  $i^{th}$ ) of  $Z_n = h(\theta_n) + \omega_n$  at a time step, and we use that to update  $i^{th}$  component of our parameter estimate:

$$\theta_{n+1}[i] = \theta_n[i] + \alpha_n Z_n[i]$$

- The convergence for this method will be proven similarly to the synchronous version, under the assumption that every coordinate is sampled infinitely often.

# Outline

- 1 TD-learning
  - TD(0)-learning
  - TD( $\lambda$ )
- 2  $Q$ -learning (tabular)
  - The  $Q$ -learning method
  - Stochastic Approximation method
- 3  $Q$ -learning with function approximation
- 4 Deep Q-learning Networks (DQN)
- 5 Approximate dynamic programming
  - TD(0) and TD(1)
  - Fitted Value iteration

# $Q$ -learning with function approximation

- The tabular  $Q$ -learning does not scale with increase in the size of state space. In most real applications, there are too many states to keep visit, and keep track of.
- For scalability, we want to **generalize**, i.e., use what we have learned about already visited (relatively small number of) states, and generalize it to new, similar states.
- A fundamental idea is to use '**function approximation**', i.e., use a lower dimensional feature representation of the state- action pair  $s, a$  and learn a parametric approximation  $Q_{\theta}(s, a)$ .

# $Q$ -learning with function approximation

- For example, the function  $Q_\theta(s, a)$  can simply be a linear function in  $\theta$  and features  $Q_\theta(s, a) = \theta_0 f_0(s, a) + \theta_1 f_1(s, a) + \dots + \theta_n f_n(s, a)$ , or a deep neural net. Given parameter  $\theta$ , the  $Q$ -function can be computed for unseen  $s, a$ . Instead of learning the  $|S| \times |A|$  dimensional  $Q$ -table, the  $Q$ -learning algorithm will learn the parameter  $\theta$ . Here, on observing sample transition to  $s'$  from  $s$  on playing action  $a$ , instead of updating the estimate of  $Q(s, a)$  in the  $Q$ -table, the algorithm updates the estimate of  $\theta$ .
- Intuitively, we are trying to find a  $\theta$  such that for every  $s, a$  the Bellman equation,

$$Q_\theta(s, a) = \mathbb{E}_{s' \sim P(\cdot|s,a)} \left[ R(s, a, s') + \gamma \max_{a'} Q_\theta(s', a') \right]$$

can be approximated well for all  $s, a$ .

# Q-learning with function approximation

- Or, in other words, we are trying to minimize some loss function like squares loss:

$$\begin{aligned}\ell_{\theta}(s, a) &= \mathbb{E}_{s' \sim P(\cdot | s, a)} \left( Q_{\theta}(s, a) - R(s, a, s') - \gamma \max_{a'} Q_{\theta}(s', a') \right)^2 \\ &= \mathbb{E}_{s' \sim P(s, a, \cdot)} [\ell_{\theta}(s, a, s')]\end{aligned}$$

where

$$\begin{aligned}\ell_{\theta}(s, a, s') &= (Q_{\theta}(s, a) - \text{target}(s'))^2 \\ \text{target}(s') &= R(s, a, s') + \gamma \max_{a'} Q_{\theta}(s', a')\end{aligned}$$

- Fundamentally, Q-learning uses gradient descent to optimize this loss function given sample (target) observations at  $s'$ .

# Q-learning Algorithm overview

Start with initial state  $s = s_0$ . In iteration  $k = 1, 2, \dots$ ,

- Take an action  $a$ .
- Observe reward  $r$ , transition to state  $s' \sim P(\cdot|s, a)$ .
- $\theta_{k+1} \leftarrow \theta_k + \alpha_k \nabla_{\theta_k} \ell_{\theta_k}(s, a, s')$ , where

$$\nabla_{\theta} \ell_{\theta_k}(s, a, s') = \delta \nabla_{\theta_k} Q_{\theta_k}(s, a)$$

$$\delta = r + \gamma \max_{a'} Q_{\theta_k}(s', a') - Q_{\theta_k}(s, a)$$

- $s \leftarrow s'$ ,

If  $s'$  reached at some point is a terminal state,  $s$  is reset to starting state.

# Outline

- 1 TD-learning
  - TD(0)-learning
  - TD( $\lambda$ )
- 2  $Q$ -learning (tabular)
  - The  $Q$ -learning method
  - Stochastic Approximation method
- 3  $Q$ -learning with function approximation
- 4 **Deep Q-learning Networks (DQN)**
- 5 Approximate dynamic programming
  - TD(0) and TD(1)
  - Fitted Value iteration

# Deep Q-Networks — Algorithm

---

**Algorithm 1** Deep Q-learning with Experience Replay

---

Initialize replay memory  $\mathcal{D}$  to capacity  $N$

Initialize action-value function  $Q$  with random weights

**for** episode = 1,  $M$  **do**

    Initialise sequence  $s_1 = \{x_1\}$  and preprocessed sequenced  $\phi_1 = \phi(s_1)$

**for**  $t = 1, T$  **do**

        With probability  $\epsilon$  select a random action  $a_t$

        otherwise select  $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$

        Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$

        Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$

        Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$

        Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to equation [3]

**end for**

**end for**

---



# Deep Deterministic Policy Gradient

- Deep Deterministic Policy Gradient (DDPG) is an algorithm which concurrently learns a Q-function and a policy. It uses off-policy data and the Bellman equation to learn the Q-function, and uses the Q-function to learn the policy.
- This approach is closely connected to Q-learning, and is motivated the same way: if you know the optimal action-value function  $Q^*(s, a)$ , then in any given state, the optimal action  $a^*(s)$  can be found by solving

$$a^*(s) = \arg \max_a Q^*(s, a).$$

- DDPG interleaves learning an approximator to  $Q^*(s, a)$  with learning an approximator to  $a^*(s)$ , and it does so in a way which is specifically adapted for environments with continuous action spaces. But what does it mean that DDPG is adapted \*specifically\* for environments with continuous action spaces? It relates to how we compute the max over actions in  $\max_a Q^*(s, a)$ .

- When there are a finite number of discrete actions, the max poses no problem, because we can just compute the Q-values for each action separately and directly compare them. (This also immediately gives us the action which maximizes the Q-value.) But when the action space is continuous, we can't exhaustively evaluate the space, and solving the optimization problem is highly non-trivial. Using a normal optimization algorithm would make calculating  $\max_a Q^*(s, a)$  a painfully expensive subroutine. And since it would need to be run every time the agent wants to take an action in the environment, this is unacceptable.
- Because the action space is continuous, the function  $Q^*(s, a)$  is presumed to be differentiable with respect to the action argument. This allows us to set up an efficient, gradient-based learning rule for a policy  $\mu(s)$  which exploits that fact. Then, instead of running an expensive optimization subroutine each time we wish to compute  $\max_a Q(s, a)$ , we can approximate it with  $\max_a Q(s, a) \approx Q(s, \mu(s))$ .

# The Q-Learning Side of DDPG

- First, let's recap the Bellman equation describing the optimal action-value function,  $Q^*(s, a)$ . It's given by

$$Q^*(s, a) = \mathbb{E}_{s' \sim P} \left[ r(s, a) + \gamma \max_{a'} Q^*(s', a') \right]$$

where  $s' \sim P$  is shorthand for saying that the next state,  $s'$ , is sampled by the environment from a distribution  $P(\cdot|s, a)$ .

- This Bellman equation is the starting point for learning an approximator to  $Q^*(s, a)$ . Suppose the approximator is a neural network  $Q_\phi(s, a)$ , with parameters  $\phi$ , and that we have collected a set  $\mathcal{D}$  of transitions  $(s, a, r, s', d)$  (where  $d$  indicates whether state  $s'$  is terminal). We can set up a **mean-squared Bellman error (MSBE)** function, which tells us roughly how closely  $Q_\phi$  comes to satisfying the Bellman equation:

$$L(\phi, \mathcal{D}) = \mathbb{E}_{(s, a, r, s', d) \sim \mathcal{D}} \left[ \left( Q_\phi(s, a) - \left( r + \gamma(1 - d) \max_{a'} Q_\phi(s', a') \right) \right)^2 \right]$$

- Here, in evaluating  $(1 - d)$ , we've used: "True" to 1 and "False" to zero. Thus, when " $d=True$ "—which is to say, when  $s'$  is a terminal state—the Q-function should show that the agent gets no additional rewards after the current state.
- Q-learning algorithms for function approximators, such as DQN (and all its variants) and DDPG, are largely based on minimizing this MSBE loss function. There are two main tricks employed by all of them which are worth describing, and then a specific detail for DDPG.
- **Trick One: Replay Buffers.** All standard algorithms for training a deep neural network to approximate  $Q^*(s, a)$  make use of an experience replay buffer. This is the set  $\mathcal{D}$  of previous experiences. In order for the algorithm to have stable behavior, the replay buffer should be large enough to contain a wide range of experiences, but it may not always be good to keep everything. If you only use the very-most recent data, you will overfit to that and things will break; if you use too much experience, you may slow down your learning. This may take some tuning to get right.

- **Trick Two: Target Networks.** Q-learning algorithms make use of **target networks**. The term

$$r + \gamma(1 - d) \max_{a'} Q_{\phi}(s', a')$$

is called the **target**, because when we minimize the MSBE loss, we are trying to make the Q-function be more like this target. Problematically, the target depends on the same parameters we are trying to train:  $\phi$ . This makes MSBE minimization unstable. The solution is to use a set of parameters which comes close to  $\phi$ , but with a time delay—that is to say, a second network, called the target network, which lags the first. The parameters of the target network are denoted  $\phi_{\text{targ}}$ .

- In DQN-based algorithms, the target network is just copied over from the main network every some-fixed-number of steps. In DDPG-style algorithms, the target network is updated once per main network update by polyak averaging:

$$\phi_{\text{targ}} \leftarrow \rho \phi_{\text{targ}} + (1 - \rho) \phi$$

where  $\rho$  is between 0 and 1 (usually close to 1).

- **DDPG Detail: Calculating the Max Over Actions in the Target.** As mentioned earlier: computing the maximum over actions in the target is a challenge in continuous action spaces. DDPG deals with this by using a **target policy network** to compute an action which approximately maximizes  $Q_{\phi_{\text{targ}}}$ . The target policy network is found the same way as the target Q-function: by polyak averaging the policy parameters over the course of training.

- Putting it all together, Q-learning in DDPG is performed by minimizing the following MSBE loss with stochastic gradient descent:

$$L(\phi, \mathcal{D}) = \mathbb{E}_{(s,a,r,s',d) \sim \mathcal{D}} \left[ \left( Q_{\phi}(s,a) - (r + \gamma(1-d)Q_{\phi_{\text{targ}}}(s', \mu_{\theta_{\text{targ}}}(s')))) \right)^2 \right],$$

where  $\mu_{\theta_{\text{targ}}}$  is the target policy.

# The Policy Learning Side of DDPG

- Policy learning in DDPG is fairly simple. We want to learn a deterministic policy  $\mu_\theta(s)$  which gives the action that maximizes  $Q_\phi(s, a)$ . Because the action space is continuous, and we assume the Q-function is differentiable with respect to action, we can just perform gradient ascent (with respect to policy parameters only) to solve

$$\max_{\theta} \mathbb{E}_{s \sim \mathcal{D}} [Q_\phi(s, \mu_\theta(s))] .$$

- Note that the Q-function parameters are treated as constants here.



# Pseudocode: DDPG

## Algorithm 3 Deep Deterministic Policy Gradient

```
1: Input: initial policy parameters  $\theta$ , Q-function parameters  $\phi$ , empty replay buffer  $\mathcal{D}$ 
2: Set target parameters equal to main parameters  $\theta_{\text{targ}} \leftarrow \theta$ ,  $\phi_{\text{targ}} \leftarrow \phi$ 
3: repeat
4:   Observe state  $s$  and select action  $a = \text{clip}(\mu_{\theta}(s) + \epsilon, a_{\text{Low}}, a_{\text{High}})$ , where  $\epsilon \sim \mathcal{N}$ 
5:   Execute  $a$  in the environment
6:   Observe next state  $s'$ , reward  $r$ , and done signal  $d$  to indicate whether  $s'$  is terminal
7:   Store  $(s, a, r, s', d)$  in replay buffer  $\mathcal{D}$ 
8:   If  $s'$  is terminal, reset environment state.
9:   if it's time to update then
10:     for however many updates do
11:       Randomly sample a batch of transitions,  $B = \{(s, a, r, s', d)\}$  from  $\mathcal{D}$ 
12:       Compute targets  $y(r, s', d) = r + \gamma(1 - d)Q_{\phi_{\text{targ}}}(s', \mu_{\theta_{\text{targ}}}(s'))$ 
13:       Update Q-function by one step of gradient descent using  $\nabla_{\phi} \frac{1}{|B|} \sum_{(s,a,r,s',d) \in B} (Q_{\phi}(s, a) - y(r, s', d))^2$ 
14:       Update policy by one step of gradient ascent using  $\nabla_{\theta} \frac{1}{|B|} \sum_{s \in B} Q_{\phi}(s, \mu_{\theta}(s))$ 
15:       Update target networks with
         
$$\phi_{\text{targ}} \leftarrow \rho \phi_{\text{targ}} + (1 - \rho) \phi$$

         
$$\theta_{\text{targ}} \leftarrow \rho \theta_{\text{targ}} + (1 - \rho) \theta$$

16:     end for
17:   end if
18: until convergence
```

# Twin Delayed DDPG (TD3)

- A common failure mode for DDPG is that the learned Q-function begins to dramatically overestimate Q-values, which then leads to the policy breaking, because it exploits the errors in the Q-function.
- **Trick One: Clipped Double-Q Learning.** TD3 learns *two* Q-functions instead of one (hence "twin"), and uses the smaller of the two Q-values to form the targets in the Bellman error loss functions.
- **Trick Two: "Delayed" Policy Updates.** TD3 updates the policy (and target networks) less frequently than the Q-function. The paper recommends one policy update for every two Q-function updates.
- **Trick Three: Target Policy Smoothing.** TD3 adds noise to the target action, to make it harder for the policy to exploit Q-function errors by smoothing out Q along changes in action.

# Key Equations: target policy smoothing

TD3 concurrently learns two Q-functions,  $Q_{\phi_1}$  and  $Q_{\phi_2}$ , by mean square Bellman error minimization, in almost the same way that DDPG learns its single Q-function.

- **target policy smoothing.** Actions used to form the Q-learning target are based on the target policy,  $\mu_{\theta_{\text{targ}}}$ , but with clipped noise added on each dimension of the action. After adding the clipped noise, the target action is then clipped to lie in the valid action range (all valid actions,  $a$ , satisfy  $a_{\text{Low}} \leq a \leq a_{\text{High}}$ ). The target actions are thus:

$$a'(s') = \text{clip}(\mu_{\theta_{\text{targ}}}(s') + \text{clip}(\epsilon, -c, c), a_{\text{Low}}, a_{\text{High}}), \quad \epsilon \sim \mathcal{N}(0, \sigma)$$

- Target policy smoothing essentially serves as a regularizer for the algorithm. It addresses a particular failure mode that can happen in DDPG: if the Q-function approximator develops an incorrect sharp peak for some actions, the policy will quickly exploit that peak and then have brittle or incorrect behavior. This can be averted by smoothing out the Q-function over similar actions, which target policy smoothing is designed to do.

# clipped double-Q learning

- Both Q-functions use a single target, calculated using whichever of the two Q-functions gives a smaller target value:

$$y(r, s', d) = r + \gamma(1 - d) \min_{i=1,2} Q_{\phi_i, \text{targ}}(s', a'(s')),$$

and then both are learned by regressing to this target:

$$L(\phi_1, \mathcal{D}) = \mathbf{E}_{(s,a,r,s',d) \sim \mathcal{D}} \left[ \left( Q_{\phi_1}(s, a) - y(r, s', d) \right)^2 \right],$$
$$L(\phi_2, \mathcal{D}) = \mathbf{E}_{(s,a,r,s',d) \sim \mathcal{D}} \left[ \left( Q_{\phi_2}(s, a) - y(r, s', d) \right)^2 \right].$$

- Using the smaller Q-value for the target, and regressing towards that, helps fend off overestimation in the Q-function.

- the policy is learned just by maximizing  $Q_{\phi_1}$ :

$$\max_{\theta} \mathbb{E}_{s \sim \mathcal{D}} [Q_{\phi_1}(s, \mu_{\theta}(s))],$$

which is pretty much unchanged from DDPG. However, in TD3, the policy is updated less frequently than the Q-functions are. This helps damp the volatility that normally arises in DDPG because of how a policy update changes the target.

- **Exploration vs. Exploitation:** TD3 trains a deterministic policy in an off-policy way. Because the policy is deterministic, if the agent were to explore on-policy, in the beginning it would probably not try a wide enough variety of actions to find useful learning signals. To make TD3 policies explore better, we add noise to their actions at training time, typically uncorrelated mean-zero Gaussian noise. To facilitate getting higher-quality training data, you may reduce the scale of the noise over the course of training.

## Algorithm 4 Twin Delayed DDPG

```
1: Input: initial policy  $\theta$ , Q-function  $\phi_1, \phi_2$ , empty replay buffer  $\mathcal{D}$ . Set  $\theta_{\text{targ}} \leftarrow \theta, \phi_{\text{targ},1} \leftarrow \phi_1, \phi_{\text{targ},2} \leftarrow \phi_2$ 
2: repeat
3:   Observe state  $s$  and select action  $a = \text{clip}(\mu_{\theta}(s) + \epsilon, a_{\text{Low}}, a_{\text{High}})$ , where  $\epsilon \sim \mathcal{N}$ 
4:   Execute  $a$  in the environment
5:   Observe next state  $s'$ , reward  $r$ , and done signal  $d$  to indicate whether  $s'$  is terminal
6:   Store  $(s, a, r, s', d)$  in replay buffer  $\mathcal{D}$ 
7:   If  $s'$  is terminal, reset environment state.
8:   if it's time to update then
9:     for  $j$  in range(however many updates) do
10:      Randomly sample a batch of transitions,  $B = \{(s, a, r, s', d)\}$  from  $\mathcal{D}$ 
11:      Compute target actions  $a'(s') = \text{clip}(\mu_{\theta_{\text{targ}}}(s') + \text{clip}(\epsilon, -c, c), a_{\text{Low}}, a_{\text{High}})$ ,  $\epsilon \sim \mathcal{N}(0, \sigma)$ 
12:      Compute targets  $y(r, s', d) = r + \gamma(1 - d) \min_{i=1,2} Q_{\phi_{\text{targ},i}}(s', a'(s'))$ 
13:      Update Q-functions by one gradient step:  $\nabla_{\phi_i} \frac{1}{|B|} \sum_{(s,a,r,s',d) \in B} (Q_{\phi,i}(s, a) - y(r, s', d))^2$  for  $i = 1, 2$ 
14:      if  $j \bmod \text{policy\_delay} = 0$  then
15:        Update policy by one step of gradient ascent using  $\nabla_{\theta} \frac{1}{|B|} \sum_{s \in B} Q_{\phi,1}(s, \mu_{\theta}(s))$ 
16:        Update target networks with
          
$$\phi_{\text{targ},i} \leftarrow \rho \phi_{\text{targ},i} + (1 - \rho) \phi_i \text{ for } i = 1, 2, \quad \theta_{\text{targ}} \leftarrow \rho \theta_{\text{targ}} + (1 - \rho) \theta$$

17:      end if
18:    end for
19:  end if
20: until convergence
```

# Exploration

- Note that we are trying to minimize many loss functions (one for each  $s, a$ ) simultaneously. At any point, the number of samples available for different state and action pairs are different, and depend on how much we explore a given state and that action. Therefore, depending on the both the transition dynamics and our exploration scheme, at any point, we will have different level of accuracies for different state and actions.
- Note that not all  $s, a$  need to have high accuracy in order to find optimal policy. For example, if some states are rare and have low value, in those states loss function does not need to be optimized. Further, if two states/action pairs have similar features, then both don't need to be explored in order to learn a good parameter  $\theta$ . This points to significance of an adaptive exploration scheme which manages appropriate accuracy levels across state and actions in order to quickly converge to a good  $\theta$ .

# Exploration

- A simple exploration scheme often utilized is  $\epsilon$ -greedy strategy, possibly with decreasing  $\epsilon$  to adapt to less exploration requirement in the later part of execution. We will study more advanced adaptive exploration methods later in the course.



# Stabilizing

- When comparing to supervised learning, we may view each sample  $s'$  as a training data row, with the loss in  $s, a$  for parameter  $\theta$  for this row being  $(Q_\theta(s, a) - \text{target}(s'))^2$ .
- However, there are several challenges compared to supervised learning. Firstly, this is different from minimizing loss compared to 'true' labels in the training data, as in supervised learning. Here, the target  $\text{target}(s') = \mathbb{E}_{s' \sim P(s'|s,a)} [R(s, a, s') + \gamma \max_{a'} Q_\theta(s', a')]$  that plays the role of labels is in fact also an estimate. This can make the learning unstable, as the target label may change too quickly as  $\theta$  changes.

# Stabilizing

- There are two main ideas utilized to make the  $Q$ -learning stable.
  - Batch learning (experience replay): Experience Replay (introduced in Lin [1992]) stores experiences including state transitions, rewards and actions, the sample observations the necessary data to perform  $Q$  learning updates. Then, these experiences are used in mini-batches to update neural networks. This increases speed of update due to batch updates, as well reduces correlation between samples used to update the parameters with decisions made from those updates.
  - Lazy update of target network: The target function is changed frequently with DNN update. Unstable target function makes training difficult. So lazy update fixes the parameters of target function and replaces them with the latest network occasionally, every thousands steps.

# Outline

- 1 TD-learning
  - TD(0)-learning
  - TD( $\lambda$ )
- 2  $Q$ -learning (tabular)
  - The  $Q$ -learning method
  - Stochastic Approximation method
- 3  $Q$ -learning with function approximation
- 4 Deep Q-learning Networks (DQN)
- 5 Approximate dynamic programming
  - TD(0) and TD(1)
  - Fitted Value iteration

## TD(0) and TD(1)

- In TD-learning we are only interested in evaluating a policy given  $\pi$  - you can think of this a special case of  $Q$ -learning on an MDP with only one action available for every state. Therefore, we are interested in computing value function  $V^\pi(s)$  for all  $s$ .
- In the function approximation version, we learn a parametric approximation  $\tilde{V}_\theta(s)$ . For example, the function  $\tilde{V}_\theta(s, a)$  could simply be a linear function in  $\theta$  and features  $\tilde{V}_\theta(s) = \theta_0 f_0(s) + \theta_1 f_1(s) + \dots + \theta_n f_n(s)$ , or output of a deep neural net. The parameter  $\theta$  can map the feature vector  $f(s)$  for any  $s$  to its value-function, thus providing a compact representation. Instead of learning the  $|S|$  dimensional value vector TD-learning algorithm will learn the parameter  $\theta$ .
- Here is a least squares approximation based modification (similar to  $Q$ -learning with function approximation) of TD(0) and TD(1) respectively.

# Approximate TD(0) method for policy evaluation

---

## Algorithm 5 Approximate TD(0) method for policy evaluation

---

- 1: Initialization: Given a starting state distribute on  $D_0$ , policy  $\pi$ , the method evaluates  $V^\pi(s)$  for all states  $s$ . Initialize  $\theta$ .
  - 2: **repeat**
  - 3:   Set  $t = 1$ ,  $s_1 \sim D_0$ . Choose step sizes  $\alpha_1, \alpha_2, \dots$ .
  - 4:   Perform TD(0) updates over an episode:
  - 5:   **repeat**
  - 6:     Take action at  $a_t \sim \pi(s_t)$ . Observe reward  $r_t$ , and new state  $s_{t+1}$ .
  - 7:     **Let**  $\delta_t := r_t + \gamma \hat{V}_\theta(s_{t+1}) - \tilde{V}_\theta(s_t)$
  - 8:     **Update**  $\theta \leftarrow \theta + \alpha_t \delta_t \nabla_\theta \tilde{V}_\theta(s_t)$
  - 9:      $t = t + 1$
  - 10:   **until** episode terminates
  - 11: **until** termination condition
-

# Fitted Value iteration

- Suppose we know the MDP model or can simulate the MDP model in any state and action. If state space was small, we could use value iteration for computing optimal policy.

## Tabular value-iteration (Recall)

- Start with an arbitrary initialization  $v^0$ . Specify  $\epsilon > 0$
- **Repeat** for  $k = 1, 2, \dots$  **until**  $\|\mathbf{v}^k(s) - \mathbf{v}^{k-1}(s)\|_\infty \leq \epsilon \frac{(1-\gamma)}{2\gamma}$ :
  - for every  $s \in S$ , improve the value vector as:

$$\mathbf{v}^k(s) = [Lv^{k-1}](s) := \max_{a \in A} R(s, a) + \gamma \sum_{s'} P(s, a, s') v^{k-1}(s')$$

# Fitted value-iteration

However, to handle large state space, we consider approximate value iteration, also known as 'fitted' value iteration.

We replace this to a scalable update, by fitting a compact representation  $\tilde{V}_\theta$  to  $v_i$  in every iteration.

## Fitted value-iteration

- Start with an arbitrary initialization  $\theta^0$ .
- **Repeat** for  $k = 1, 2, 3, \dots$ :
  - Evaluate  $[L\tilde{V}_{\theta^{k-1}}](s)$  on a subset  $s \in S_0$ .
  - Use some regression technique to find a  $\theta$  to fit data  $(\tilde{V}_\theta(s), L\tilde{V}_{\theta^{k-1}}(s))_{s \in S_0}$ . Set this  $\theta$  as  $\theta_k$ .

## Example 1: Sampling + Least squares fitting

Sampling + Least squares fitting [Munos and Szepesvari, 2008] (This assumes number of actions is small, number of states is large. And, the MDP model is large but can be simulated for any arbitrary state  $s$  and action  $a$ )

In state  $k$ :

- Sample states  $X_1, X_2, \dots, X_N$  from the state space  $S$ , using some distribution  $\mu$ .
- For each action  $a$ , and state  $X_i$ , take multiple samples of next state and rewards  $\{Y_{i,a,j}, R_{i,a,j}\}_{j=1, \dots, M}$ .
- Approximate  $[L\tilde{V}_{\theta^{k-1}}](s)$  for  $s \in \{X_1, X_2, \dots, X_N\}$  as

$$\text{target}_i = \max_{a \in A} \frac{1}{M} \sum_{j=1}^M (R_{i,a,j} + \gamma \tilde{V}_{\theta^{k-1}}(Y_{i,a,j})), \forall i = 1, \dots, N$$

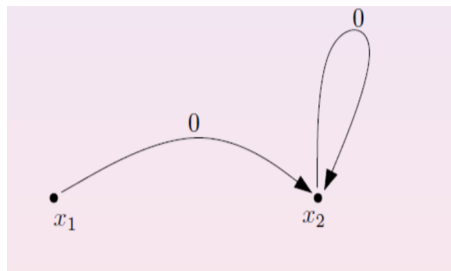
- Use least squares new  $\theta$  as:

$$\theta_k := \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N (\tilde{V}_{\theta}(i) - \text{target}_i)^2$$



# Counter-example for least-square regression

Counter-example for least-square regression [Tsitsiklis and van Roy, 1996] An MDP with two states  $x_1, x_2$ , 1-d features for the two states:  $f_{x_1} = 1, f_{x_2} = 2$ . Linear Function approximation with  $\tilde{V}_\theta(x) = \theta f_x$ .



# Counter-example for least-square regression

$$\begin{aligned}\theta_k : &= \arg \min_{\theta} \frac{1}{2} (\theta - \text{target}_1)^2 + (2\theta - \text{target}_2)^2 \\ &= \arg \min_{\theta} \frac{1}{2} (\theta - \gamma \theta^{k-1} f_{x_2})^2 + (2\theta - \gamma \theta^{k-1} f_{x_2})^2 \\ &= \arg \min_{\theta} \frac{1}{2} (\theta - \gamma 2\theta^{k-1})^2 + (2\theta - \gamma 2\theta^{k-1})^2\end{aligned}$$

$$(\theta - \gamma 2\theta^{k-1}) + 2(2\theta - \gamma 2\theta^{k-1}) = 0 \Rightarrow 5\theta = 6\gamma \theta^{k-1}$$

$$\theta_k = \frac{6}{5} \gamma \theta_{k-1}$$

This diverges if  $\gamma \geq 5/6$ .

# Operator view of Fitted value-iteration

Consider operator  $M_A$  defined as follows: Fit a  $\tilde{V}_\theta$  to  $L\tilde{V}_{\theta^{i-1}}$  by comparing its values on a subset  $S_0$  of states, using a regression technique. Then return this  $\tilde{V}_\theta$  as new function  $\tilde{V}_{\theta^i}$  in the output space of  $M_A$ . Thus,  $M_A$  is effectively an approximation operator.

Equivalently,

- Start with an arbitrary initialization  $v^0$ .
- **Repeat** for  $k = 1, 2, 3, \dots$ :
  - $v^i = (L \circ M_A) v^{i-1}$ .

(In an efficient implementation,  $u^{i-1} = M_A v^{i-1}$  probably has a more compact representation, so the first view may be better for implementation)

# Operator view of Fitted value-iteration

The above view allows us to view fitted value iteration as just value iteration with a different operator:  $v^i = Lv^{i-1}$  is replaced by  $\tilde{V}^i = (M_A \circ L) \tilde{V}^{i-1}$ . Therefore, as long as the new operator  $(M_A \circ L)$  is also  $\gamma$ -contraction, the results proven earlier for convergence of value iteration will hold. A sufficient condition is that the operator  $M_A$  is non-expansive:

$$\|M_A v - M_A v'\|_\infty \leq \|v - v'\|_\infty$$

Then,

$$\begin{aligned} \|(M_A \circ L) v - (M_A \circ L) v'\|_\infty &\leq \|Lv - Lv'\|_\infty \\ &\leq \gamma \|v - v'\|_\infty \end{aligned}$$

Similarly, in the other view,  $v^i = Lv^{i-1}$  is replaced by  $v^i = (T \circ M_A) v^{i-1}$  which is also a  $\gamma$  contraction. So,  $v^i$  converges.

# Sufficient condition for convergence

Sufficient condition for convergence: Averager [Gordon, 1995].

## Lemma 1

The operator  $M_A$  is non-expansive if it is an averager, that is,

$$[M_A v^k](s) = \sum_{i \in S} w_{i,s} v^k(i) + w_{0,s}$$

where,

$$\sum_i w_{i,s} + w_{0,s} = 1, w_{i,s} \geq 0$$

# Proof of Convergence

**Proof.** This is non-expansive because

$$M_A v = Wv + w_0$$

$$\begin{aligned}\|M_A v - M_A v'\|_\infty &= \max_s \left| \sum_{i \in S} w_{i,s} (v(i) - v'(i)) \right| \\ &\leq \max_s \max_i |v(i) - v'(i)| \\ &= \|v - v'\|_\infty\end{aligned}$$

## Converges to what

The fitted value iteration converges with exponential rate under above condition, but to what? Let  $v^i$  converges to fixed point  $U^*$ . That is,  $U^* = L(M_A(U^*))$ . We show that under the contraction properties,  $U^*$  is as close as it can be to  $V^*$  (the optimal value function), in the sense that:

$$\|U^* - V^*\|_\infty \leq \frac{\gamma}{1 - \gamma} \|M_A V^* - V^*\|_\infty$$

The proof is as follows.

$$\begin{aligned}\|U^* - V^*\| &= \|L(M_A(U^*)) - LV^*\| \\ &\leq \gamma \|M_A U^* - V^*\| \\ \|M_A U^* - V^*\| &\leq \|M_A U^* - M_A V^*\| + \|M_A V^* - V^*\| \\ &\leq \|U^* - V^*\| + \|M_A V^* - V^*\| \\ \|U^* - V^*\| &\leq \gamma \|U^* - V^*\| + \gamma \|M_A V^* - V^*\| \\ \|U^* - V^*\| &\leq \frac{\gamma}{(1 - \gamma)} \|M_A V^* - V^*\|\end{aligned}$$