



**RLChina 2020**

# Lecture 2: Foundations of Reinforcement Learning

— from policy methods to PAC bounds analysis

Prof. Jun Wang  
UCL

July 28, 2020

# Table of content

- ▶ Recap (yesterday's lecture)
  - ▶ MDPs
  - ▶ Value iterations and policy iterations
  - ▶ Tabular Q-Learning
- ▶ Policy approaches
  - ▶ Markov chains
  - ▶ Policy gradient
- ▶ Computational learning theory
  - ▶ PAC learning concepts
  - ▶ Learning bound for finite  $H$
- ▶ Theoretical analysis
  - ▶ Approximate dynamic programming
  - ▶ Performance bounds
  - ▶ Sample complexity

# Main references

- ▶ machine learning and learning theory books<sup>12</sup>
- ▶ reinforcement learning books<sup>34</sup>
- ▶ approximate dynamic programming<sup>45</sup>
- ▶ this slide is adopted from our upcoming book chapter<sup>6</sup>

---


<sup>1</sup>Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.

<sup>2</sup>Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.

<sup>3</sup>Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

<sup>4</sup>Dimitri P Bertsekas and John N Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, 1996.

<sup>5</sup>Rémi Munos. *Introduction to Reinforcement Learning and multi-armed bandits*. NETADIS Summer School, 2013.

<sup>6</sup>Shuang Wu and Jun Wang. *Decision making and AI: a white paper*. 2020.  3/73

# Table of content

- ▶ **Recap**
  - ▶ MDPs
  - ▶ Value iterations and policy iterations
  - ▶ Tabular Q-Learning
- ▶ Policy approaches
  - ▶ Markov chains
  - ▶ Policy gradient
- ▶ Computational learning theory
  - ▶ PAC learning concepts
  - ▶ Learning bound for finite  $H$
- ▶ Theoretical analysis
  - ▶ Approximate dynamic programming
  - ▶ Performance bounds
  - ▶ Sample complexity

# Markov Decision Process(MDP)

## Definition:

- ▶ an MDP  $M$  is a tuple  $\{\mathbb{S}, \mathbb{A}, \Pr(s'|s, a), r(s, a, s')\}$
- ▶  $\mathbb{S}$ : state space; a state  $s \in \mathbb{S}$
- ▶  $\mathbb{A}$ : action space; an action  $a \in \mathbb{A}$
- ▶ system dynamics:  $\Pr(s'|s, a)$
- ▶ reward:  $r(s, a, s')$ <sup>7</sup>

---

<sup>7</sup>For the MDP with known state transition  $\Pr(s'|s, a)$ , the stochastic reward  $r(s, a, s')$  can be reduced to a deterministic one as

$$R(s, a) := \mathbb{E}_{s' \sim \Pr(s'|s, a)}[r(s, a, s')].$$

# Policy

- ▶ a policy specifies what an agent should do in a specific circumstance
- ▶ it is a mapping from the history to a (deterministically or randomly) action at present:  $\pi_k : (s_0, a_0, s_1, a_1, \dots, s_k) \mapsto a_k$ 
  - ▶ *Markovian policy*  $\pi_k : s_k \mapsto a_k$
  - ▶ *stationary policy*  $\pi_k = \pi_{k+1}$  for any  $k$
  - ▶ *deterministic policy*:  $a := \pi(s)$
  - ▶ *randomized policy*  $\Pr(a|s) := \pi(s|a)$  is often considered for an RL setup as it helps explore the environments

# Objectives

an agent chooses a policy in order to maximise one of the following possible objectives:

1. total reward MDP over a finite horizon,

$$J(\pi) := \mathbb{E} \left[ \sum_{k=0}^N R(s_k, \pi(s_k)) \right];$$

2. discounted reward MDP over an infinite horizon,

$$J(\pi) := \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k R(s_k, \pi(s_k)) \right], \text{ for } 0 < \gamma < 1;$$

3. average reward MDP over an infinite horizon (ergodic reward),

$$J(\pi) = \lim_{K \rightarrow \infty} \mathbb{E} \left[ \frac{1}{K+1} \sum_{k=0}^K R(s_k, \pi(s_k)) \right]$$

# Policy evaluation and optimisation

given the objective, the two core questions in an MDP are:

1. how good is a policy?  $\implies$  **policy evaluation**
2. what is the optimal policy?  $\implies$  **policy optimisation**



# Value function and Q function

1. **policy evaluation** *value function*: the expected discounted rewards under a policy  $\pi$  starting from state  $s$ ,

$$V^\pi(s) := \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k R(s_k, \pi(s_k)) \mid s_0 = s \right],$$

and the corresponding action-value function (Q function), the expected reward of taking a particular action, under a policy

$$Q^\pi(s) := R(s, \pi(s)) + \mathbb{E}_{s' \sim \text{Pr}(s'|s, \pi(s))} [V^*(s')].$$

2. **policy optimisation** the value function optimising the policy gives the optimal value function,

$$V^*(s) := \max_{\pi} \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k R(s_k, \pi(s_k)) \mid s_0 = s \right].$$

By the same token, the optimal Q function is

$$Q^*(s) := \max_a R(s, a) + \mathbb{E}_{s' \sim \text{Pr}(s'|s, a)} [V^\pi(s')].$$

## Solving the value function and Q function

- ▶ **model-based approaches** when the state transition and the reward function are known and given, the MDP can be solved by dynamic programming (DP) [Bel57], e.g., value iteration or policy iteration
- ▶ **model-free approaches** if the state transition and the reward function are not given, one can solve MDP with *reinforcement learning* (RL) by *learning* the solution through interactions with the environment

## Bellman equation: how good a policy is

evaluation is also known as prediction. by splitting the immediate reward from  $V^\pi(s)$  as

$$\begin{aligned} V^\pi(s) &= \mathbb{E} \left[ \gamma^0 R(s_0, \pi(s_0)) | s_0 = s \right] + \mathbb{E} \left[ \sum_{k=1}^{\infty} \gamma^k R(s_k, \pi(s_k)) | s_0 = s \right] \\ &= R(s, \pi(s)) + \mathbb{E}_{s' \sim \Pr(s'|s, \pi(s))} \left[ \sum_{k=1}^{\infty} \gamma^k R(s_k, \pi(s_k)) | s_1 = s' \right] \\ &= R(s, \pi(s)) + \gamma \sum_{s'} \Pr(s'|s, a) \underbrace{\mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^k R(s_k, \pi(s_k)) | s_0 = s' \right]}_{=V^\pi(s')} \end{aligned}$$

we can obtain the *Bellman equation* for policy  $\pi$  as

$$V^\pi(s) = R(s, \pi(s)) + \underbrace{\gamma \sum_{s'} \Pr(s'|s, \pi(s)) V^\pi(s')}_{\text{cost-to-go}}. \quad (1)$$

## Value iteration for evaluating a policy

when  $\Pr(s'|s, a)$  and  $R(s, a)$  are given,  $V^\pi$  can be computed with either value iteration sketched in Algorithm 1

---

**Algorithm 1** value iteration for evaluating a policy

---

- 1: initialise  $\pi$  and  $V$  arbitrarily
  - 2: **repeat**
  - 3:      $V(s) \leftarrow R(s, \pi(s)) + \gamma \sum_{s'} \Pr(s'|s, \pi(s))V(s'), \forall s$
  - 4: **until** convergence
- 

both  $\pi$  and  $V$  functions can be as a form of tables (tabular).

## Bellman optimality equation: optimal policy

optimisation is also called as control. by splitting the reward, we can obtain the *Bellman optimality equation* as

$$V^*(s) = \max_{a \in \mathbb{A}} R(s, a) + \underbrace{\gamma \sum_{s'} \Pr(s'|s, a) V^*(s')}_{\text{cost-to-go}}. \quad (2)$$

when  $\Pr(s'|s, a)$  and  $R(s, a)$  are given, the optimal value function and the optimal policy can be directly solved with value iteration [Bel57] and policy iteration [How60], respectively.

$$\text{(value iteration)} \quad V(s) \leftarrow \max_{a \in \mathbb{A}} R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V(s') \quad (3)$$

$$\text{(policy iteration)} \quad \pi(s) \leftarrow \arg \max_{a \in \mathbb{A}} R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V^\pi(s') \quad (4)$$

---

**Algorithm 2** Value iteration for optimising a policy
 

---

- 1: Initialize  $\pi$  and  $V$  arbitrarily
  - 2: **repeat**
  - 3:      $V(s) \leftarrow \max_{a \in \mathbb{A}} R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V(s'), \forall s$      ▷  
       value improvement
  - 4: **until** convergence
- 

- ▶ value iteration computes the optimal value function  $V^*$  under the optimal policy
- ▶ using  $V^*(s)$ , the optimal policy  $\pi^*$  can be calculated from the Bellman optimality equation by

$$a^* = \pi^*(s) = \arg \max_{a \in \mathbb{A}} R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V^*(s').$$

---

**Algorithm 3** Policy iteration for optimising a policy

---

- 1: Initialize  $\pi$  and  $V$  arbitrarily
  - 2: **repeat**
  - 3:     **repeat**
  - 4:          $V(s) \leftarrow R(s, \pi(s)) + \gamma \sum_{s'} \Pr(s'|s, \pi(s))V(s'), \forall s$    ▶  
       Policy evaluation using Algorithm 1
  - 5:     **until** Convergence
  - 6:          $\pi(s) \leftarrow \arg \max_{a \in \mathbb{A}} R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a)V(s'), \forall s$    ▶  
       Policy improvement
  - 7: **until** convergence
- 

- ▶ policy iteration directly generates the optimal policy

## VI and PI: convergence analysis

Let  $\mathbb{F}$  be the space of functions on domain  $\mathcal{S}$ . Define the Bellman *policy* operator  $T^\pi : \mathbb{F} \mapsto \mathbb{F}$  and the Bellman *optimality* operator  $T : \mathbb{F} \mapsto \mathbb{F}$  as

$$T^\pi V(s) := R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V(s'), \quad \forall s, \quad (5)$$

$$TV(s) := \max_{a \in \mathbb{A}} R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V(s') \quad \forall s \quad (6)$$



## VI and PI: convergence analysis

- ▶ these two operators are mappings from one value function  $V(s)$  to another  $TV(s)$
- ▶ also they are both *monotonic* and *contraction* mappings (with respect to the  $\infty$ -norm), i.e.,

1. *monotonicity*, if  $V(s) \geq U(s)$  for any  $s$ ,

$$T^\pi V(s) \geq T^\pi U(s); \quad (7)$$

$$TV(s) \geq TU(s) \quad (8)$$

2. *contraction*, for any  $U, V$ ,

$$\|T^\pi V - T^\pi U\|_\infty \leq \gamma \|V - U\|_\infty \quad (9)$$

$$\|TV - TU\|_\infty \leq \gamma \|V - U\|_\infty \quad (10)$$

## Monotonicity proof

the inequality (7) follows from

$$T^\pi V(s) - T^\pi U(s) = \sum_{s'} \Pr(s'|s, \pi(s))(V(s') - U(s')) \geq 0.$$

the inequality (8) follows from, for any  $a$ ,

$$R(s, a) + \sum_{s'} \Pr(s'|s, a)V(s') \geq R(s, a) + \sum_{s'} \Pr(s'|s, a)U(s').$$

## Contraction mapping proof

The inequality (9) follows from

$$\begin{aligned}\|T^\pi V - T^\pi U\|_\infty &= \max_s \gamma \sum_{s'} \Pr(s'|s, \pi(s)) |V(s') - U(s')| \\ &\leq \gamma \left( \sum_{s'} \Pr(s'|s, \pi(s)) \right) \max_{s'} |V(s') - U(s')| \leq \gamma \|U - V\|_\infty.\end{aligned}$$

The inequality (10) follows from

$$\begin{aligned}\|TV - TU\|_\infty &= \max_s | \max_a \{R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V(s')\} \\ &\quad - \max_a \{R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) U(s')\} | \\ &\leq \max_{s,a} |R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V(s') \\ &\quad - R(s, a) - \gamma \sum_{s'} \Pr(s'|s, a) V(s')| \\ &= \gamma \max_{s,a} | \sum_{s'} \Pr(s'|s, a) (V(s') - U(s')) | \\ &\leq \gamma \left( \sum_{s'} \Pr(s'|s, a) \right) \max_{s'} |(V(s') - U(s'))| \leq \gamma \|V - U\|_\infty.\end{aligned}$$

## Contraction mapping proof

for any contraction operator (take  $T$  for example), we have a unique fixed point  $V^* = TV^*$  by Cauchy convergence theorem<sup>8</sup>. We therefore have the convergence as

$$\begin{aligned}\|V_{k+1} - V^*\|_\infty &= \|TV_k - TV^*\|_\infty \\ &\leq \gamma \|V_k - V^*\|_\infty \leq \dots \leq \gamma^{k+1} \|V_0 - V^*\|_\infty \rightarrow 0.\end{aligned}$$

---

8

- ▶ a sequence of real numbers  $(a_n)$  is said to be a Cauchy Sequence if  $\forall \epsilon > 0$  there exists an  $N \in \mathbb{N}$  such that if  $m, n \geq N$  then  $|a_n - a_m| < \epsilon$
- ▶ theorem (Cauchy Convergence Criterion): If  $(a_n)$  is a sequence of real numbers, then  $(a_n)$  is convergent if and only if  $(a_n)$  is a Cauchy sequence

## Policy iteration: convergence

- ▶ policy iteration monotonically improves the policy by

$$\begin{aligned}V^{\pi_{k+1}} &= (I - \gamma P^{\pi_{k+1}})^{-1} R^{\pi_{k+1}} \\ &\geq (I - \gamma P^{\pi_{k+1}})^{-1} (V^{\pi_k} - \gamma P^{\pi_{k+1}} V^{\pi_k}) \\ &= V^{\pi_{\text{old}}},\end{aligned}\tag{11}$$

where (11) follows from

$$\begin{aligned}R^{\pi_{k+1}} + \gamma P^{\pi_{k+1}} V^{\pi_k} &= T^{\pi_{k+1}} V^{\pi_k} = TV^{\pi_k} \geq V^{\pi_k} \\ &\iff R^{\pi_{k+1}} \geq (I - \gamma P^{\pi_{k+1}}) V^{\pi_k}.\end{aligned}$$

by the monotone convergence theorem<sup>9</sup>,  $V^{\pi^*} = \lim_{k \rightarrow \infty} V^{\pi_k}$  exists and satisfies

$$V^{\pi^*} = TV^{\pi^*},$$

which satisfies the Bellman optimality equation.

---

<sup>9</sup>If a sequence of real numbers is increasing and bounded above, then its supremum is the limit.

# Problems with VI and PI

- ▶ *continuous states and action.* both VI and PI fail when the states and actions are continuous
- ▶ *curse of dimensionality.* both VI and PI involve iterative scheme over the whole state space. The algorithm is not scalable with respect to the size of the state space
- ▶ *partial observable states.* in reality, the states may not be fully observable. the partial observability model leads to a partially observable MDP (**POMDP**)
- ▶ *unknown model.* using DP to solve an MDP requires knowledge of  $\Pr(s'|s, a)$  and  $R(s, a)$ . These quantities are costly or even impossible to acquire, especially for huge state space and action space.
  - ▶ nevertheless, the agent can collect samples of state transitions  $s, a \rightarrow s'$  and associated reward  $r(s, a, s')$  through interaction with the environment.
  - ▶ this suits machine learning and motivates the development of reinforcement learning

## Q learning

- ▶ Q-learning [Wat89] learns from the estimated optimal value function
- ▶ it defines action-value function  $Q^\pi(s, a)$  for policy  $\pi$ , which reflects the future reward for different actions under  $\pi$  as

$$\begin{aligned} Q^\pi(s, a) &= \sum_{s'} \Pr(s'|s, a) \left( r(s, a, s') + \gamma V^\pi(s') \right) \\ &= \sum_{s'} \Pr(s'|s, a) \left( r(s, a, s') + \gamma Q^\pi(s', \pi(s')) \right). \end{aligned}$$

- ▶ the optimal  $Q^*(s, a)$  is related to the optimal value function  $V^*(s) = \max_a Q^*(s, a)$ , thus

$$Q^*(s, a) = \sum_{s'} \Pr(s'|s, a) [r(s, a, s') + \gamma \max_{a'} Q(s', a')].$$

- ▶ plugging  $Q(s, a) = Q^*(s, a)$  into the value iteration for optimising policy and **replace the expectation with its sample**, we have

$$\text{(Q-learning)} \quad Q(s, a) \leftarrow Q(s, a) + \alpha \underbrace{[r + \gamma \max_{a'} Q(s', a')] - Q(s, a)}_{\text{sample of } Q^*(s, a)}.$$

# Roadmap

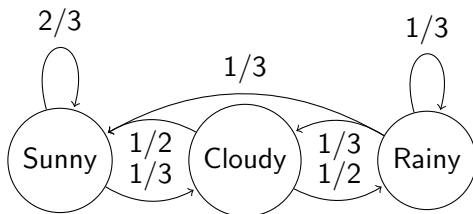


# Table of content

- ▶ Recap
  - ▶ MDPs
  - ▶ Value iterations and policy iterations
  - ▶ Tabular Q-Learning
- ▶ **Policy approaches**
  - ▶ Markov chains
  - ▶ Policy gradient
- ▶ Computational learning theory
  - ▶ PAC learning concepts
  - ▶ Learning bound for finite  $H$
- ▶ Theoretical analysis
  - ▶ Approximate dynamic programming
  - ▶ Performance bounds
  - ▶ Sample complexity

## Markov chains: an example

- ▶ *sunny*, *rainy*, and *cloudy* are called the states of the Markov chain
- ▶ if the weather is currently *sunny*, what is the prediction for the next few days according to the model?



Day	sunny	cloudy	rainy
0	1	0	0
1	$2/3$	$1/3$	0
2	0.611	0.222	0.167
3	0.574	0.259	0.167
...	0.568	0.247	0.185

## $n$ -step transition probability

- ▶ The  $n$ -step transition probability of a Markov chain is the probability that it goes from state  $i$  to state  $j$  in exactly  $n$  steps (transitions):

$$p_{ij}^{(n)} := P(S_{n+m} = j | S_m = i)$$

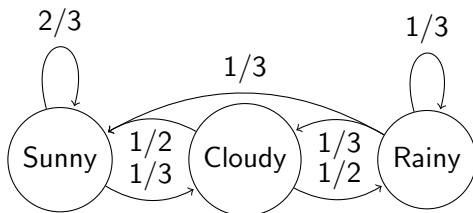
- ▶ define  $P = (p_{ij})$  the transition matrix; then the  $n$ -step transition matrix is given by  $n$  powers of the matrix:

$$P^{(n)} = P^n, \text{ for } n \geq 1$$

- ▶  $p(\mathbf{i \text{ to } j \text{ in } n \text{ steps}}) =$  sum of probabilities of all paths  $i$  to  $j$  in  $n$  steps

$$p_{ij}^{(n+m)} = \sum_k p_{ik}^{(m)} p_{kj}^{(n)}$$

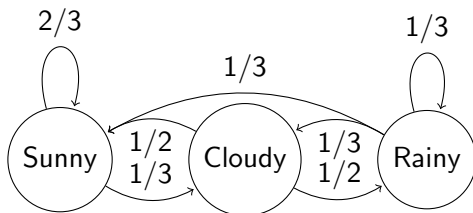
## n-step transition probability: an example



$$P^3 \approx \begin{pmatrix} 0.574 & 0.259 & 0.167 \\ 0.556 & 0.222 & 0.222 \\ 0.537 & 0.259 & 0.204 \end{pmatrix}; P^{10} \approx \begin{pmatrix} 0.563 & 0.250 & 0.187 \\ 0.562 & 0.250 & 0.187 \\ 0.562 & 0.250 & 0.188 \end{pmatrix}$$

- ▶ regardless of the initial weather  $q^{(1)}$  is,  $q^{(1)} \cdot P^n$  seems to approach  $\tilde{q} \approx (0.563, 0.250, 0.188)$  as  $n$  grows
- ▶ if we multiply the vector  $\tilde{q}$  with  $P$ , we almost get  $\tilde{q}$  again, e.g.,  $\tilde{q}$  is almost an eigenvector of  $P$  with eigenvalue 1

## n-step transition probability: an example



$$P^3 \approx \begin{pmatrix} 0.574 & 0.259 & 0.167 \\ 0.556 & 0.222 & 0.222 \\ 0.537 & 0.259 & 0.204 \end{pmatrix}; P^{10} \approx \begin{pmatrix} 0.563 & 0.250 & 0.187 \\ 0.562 & 0.250 & 0.187 \\ 0.562 & 0.250 & 0.188 \end{pmatrix}$$

- ▶ a distribution  $q$  over the states is *stationary distribution* of the Markov chain with transition matrix  $P$  if  $q = q \cdot P$
- ▶  $q$  could be interpreted as the *long-term visitation rate*



# Ergodic Markov chains

- ▶ even without dead-ends, a graph may not have well-defined long-term visit rates;
  - ▶ *requirement 1*: there is a path from any state to any other state
  - ▶ *requirement 2*: the states cannot be partitioned such that the random walker visits the partitions sequentially (no loop!)
- ▶ in an Ergodic Markov Chains:
  - ▶ The  $p^{(n)}$  has settled to a limiting value  $q$

$$p_{ij}^{(n)} \rightarrow q_j, \text{ as } n \rightarrow \infty$$

- ▶ This value is independent of initial state  $q^{(1)}$
- ▶ The  $q^{(n)}$  also approaches this limiting value:  $q_j^{(n)} \rightarrow q_j$
- ▶ where  $q$  is the unique stationary distribution of the chain (i.e. the limiting distribution is the stationary distribution)

## Go back to RL: the objective function

- ▶ suppose a policy, denoted as  $\pi_\theta$ , is parameterised by  $\theta$
- ▶ the expected reward (as the objective):

$$J(\theta) = \sum_{s \in \mathcal{S}} d^\pi(s) V^\pi(s) = \sum_{s \in \mathcal{S}} d^\pi(s) \left( \sum_{a \in \mathcal{A}} \pi_\theta(a | s) Q^\pi(s, a) \right)$$

where  $d^\pi(s) := \lim_{t \rightarrow \infty} p(S_t = s | s_0, \pi_\theta)$  is the stationary distribution of the Markov chain when the agent starts from  $s_0$  and following policy  $\pi_\theta$  for  $t$  steps

- ▶ When  $\pi_\theta$  is given, as the time progresses, the probability that the agent ends up with one state becomes unchanged (regardless of  $s_0$ ) if the underlying Markov chain is *ergodic*
- ▶ we thus drop  $s_0$  in  $J(\theta)$



# Policy gradient

the expected reward (as the objective):

$$J(\theta) = \sum_{s \in S} d^\pi(s) V^\pi(s) = \sum_{s \in S} d^\pi(s) \left( \sum_{a \in A} \pi_\theta(a | s) Q^\pi(s, a) \right)$$

where  $d^\pi(s) := \lim_{t \rightarrow \infty} p(S_t = s | s_0, \pi_\theta)$

- ▶ *gradient ascent* moves  $\theta$  toward the direction suggested by the gradient  $\nabla_\theta J(\theta)$  to find the best  $\theta$  for  $\pi_\theta$  that produces the highest return
- ▶ but, computing  $\nabla_\theta J(\theta)$  is tricky as it depends on both the action selection (directly determined by  $\pi_\theta$ ) and the stationary distribution of states  $d^\pi$  (indirectly determined by  $\pi_\theta$ )
- ▶ given that the environment is generally unknown, it is difficult to estimate the effect on the state distribution by a policy update.

# Policy gradient theorem

- ▶ **Policy gradient theorem**<sup>10</sup>: for an MDP,

$$\nabla_{\theta} J(\theta) = \sum_{s \in S} d^{\pi}(s) \left( \sum_{a \in A} Q^{\pi}(s, a) \nabla_{\theta} \pi_{\theta}(a | s) \right)$$

- ▶ It provides a nice reformation of the objective function that does not involve the derivative of the state distribution  $\frac{\partial d^{\pi}(s)}{\partial \theta}$

---

<sup>10</sup>Richard S Sutton et al. “Policy gradient methods for reinforcement learning with function approximation”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 1999, pp. 1057–1063.

# Policy gradient theorem [SB98]: the proof

$$\begin{aligned}
 & \nabla_{\theta} \left( \mathbb{E} \left[ \sum_{k=1}^{\alpha} \gamma^{k-1} r_{t+k} \mid s_t = s_0, \pi \right] \right) = \nabla_{\theta} V^{\pi}(s_0) = \nabla_{\theta} \left( \sum_{a \in A} Q^{\pi}(s_0, a) \pi_{\theta}(a \mid s_0) \right) \\
 &= \sum_{a \in A} (Q^{\pi}(s_0, a) \nabla_{\theta} \pi_{\theta}(a \mid s_0) + \pi_{\theta}(a \mid s_0) \nabla_{\theta} Q^{\pi}(s_0, a)) \quad \text{product rule} \\
 &= \sum_{a \in A} \left( Q^{\pi}(s_0, a) \nabla_{\theta} \pi_{\theta}(a \mid s_0) + \pi_{\theta}(a \mid s_0) \nabla_{\theta} \left( \sum_{s', r} P(s', r \mid s, a) (r + V^{\pi}(s')) \right) \right) \quad \text{expand} \\
 &= \sum_{a \in A} \left( Q^{\pi}(s_0, a) \nabla_{\theta} \pi_{\theta}(a \mid s_0) + \pi_{\theta}(a \mid s_0) \left( \sum_{s', r} P(s', r \mid s_0, a) \nabla_{\theta} V^{\pi}(s') \right) \right) \quad \text{remove } r \\
 &= \sum_{a \in A} \left( Q^{\pi}(s_0, a) \nabla_{\theta} \pi_{\theta}(a \mid s_0) + \pi_{\theta}(a \mid s_0) \left( \sum_{s'} P(s' \mid s_0, a) \nabla_{\theta} V^{\pi}(s') \right) \right) \quad \text{marginalise } r \text{ over } r \\
 &= \sum_{a \in A} (Q^{\pi}(s_0, a) \nabla_{\theta} \pi_{\theta}(a \mid s_0)) + \sum_{a \in A} \pi_{\theta}(a \mid s_0) \left( \sum_{s'} P(s' \mid s_0, a) \nabla_{\theta} V^{\pi}(s') \right) \\
 &= \sum_{a \in A} (Q^{\pi}(s_0, a) \nabla_{\theta} \pi_{\theta}(a \mid s_0)) + \sum_{s'} \left( \sum_{a \in A} \pi_{\theta}(a \mid s_0) P(s' \mid s_0, a) \right) \nabla_{\theta} V^{\pi}(s')
 \end{aligned}$$

## Policy gradient theorem [SB98]: the proof

- ▶ we thus have the following recursive form of the gradient:

$$\begin{aligned}\nabla_{\theta} V^{\pi}(s_0) &= \sum_{a \in A} (Q^{\pi}(s_0, a) \nabla_{\theta} \pi_{\theta}(a | s_0)) \\ &\quad + \sum_{s'} \left( \sum_{a \in A} \pi_{\theta}(a | s_0) P(s' | s_0, a) \right) \nabla_{\theta} V^{\pi}(s')\end{aligned}$$

- ▶ to simplify this, we have:

$$\nabla_{\theta} V^{\pi}(s_0) = \varphi(s_0) + \sum_{s'} P^{\pi}(s' | s_0) \nabla_{\theta} V^{\pi}(s'),$$

where  $\varphi(s_0) := \sum_{a \in A} (Q^{\pi}(s_0, a) \nabla_{\theta} \pi_{\theta}(a | s_0))$  and Markov transition  $P^{\pi}(s' | s_0) := \sum_{a \in A} \pi_{\theta}(a | s_0) P(s' | s_0, a)$

## Policy gradient theorem [SB98]: the proof

- ▶ We now consider the following visitation sequence:

$$s_0 \xrightarrow{a \sim \pi_\theta(\cdot | s_0)} s' \xrightarrow{a \sim \pi_\theta(\cdot | s')} s'' \xrightarrow{a \sim \pi_\theta(\cdot | s'')} \dots \xrightarrow{a \sim \pi_\theta(\cdot | \cdot)} s$$

- ▶ and denote the probability of transitioning from state  $s_0$  to state  $s$  with policy  $\pi_\theta$  after  $k$  step as

$$P^\pi(s'' | s_0, k) \equiv \sum_{s'} P^\pi(s'' | s', k-1) P^\pi(s' | s_0)$$

- ▶ where we have  $P^\pi(s | s_0, 0) = 1$  if  $s = s_0$  and  $P^\pi(s | s_0, 0) = 0$  if  $s \neq s_0$  (because it happened already!)
- ▶ Let us now go back to unroll the gradient of the value function:

$$\nabla_\theta V^\pi(s_0) = \varphi(s_0) + \sum_{s'} P^\pi(s' | s_0) \nabla_\theta V^\pi(s')$$

# Policy gradient theorem: the proof

let us now go back to unroll the gradient of the value function:

$$\begin{aligned}\nabla_{\theta} V^{\pi}(s_0) &= \varphi(s_0) + \sum_{s'} P^{\pi}(s' | s_0) \nabla_{\theta} V^{\pi}(s') \\ &= \varphi(s_0) + \sum_{s'} P^{\pi}(s' | s_0, 1) \left[ \varphi(s') + \sum_{s''} P^{\pi}(s'' | s') \nabla_{\theta} V^{\pi}(s'') \right] \\ &= \varphi(s_0) + \left[ \sum_{s'} P^{\pi}(s' | s_0, 1) \varphi(s') \right] + \left[ \sum_{s'} P^{\pi}(s' | s, 1) \sum_{s''} P^{\pi}(s'' | s') \nabla_{\theta} V^{\pi}(s'') \right] \\ &= \varphi(s_0) + \left[ \sum_{s'} P^{\pi}(s' | s_0, 1) \varphi(s') \right] + \left[ \sum_{s''} \sum_{s'} P^{\pi}(s' | s, 1) P^{\pi}(s'' | s') \nabla_{\theta} V^{\pi}(s'') \right] \\ &= \varphi(s_0) + \left[ \sum_{s'} P^{\pi}(s' | s_0, 1) \varphi(s') \right] + \left[ \sum_{s''} P^{\pi}(s' | s, 2) \nabla_{\theta} V^{\pi}(s'') \right] \\ &= \sum_{s \in S} \sum_{k=0}^{\infty} P^{\pi}(s | s_0, k) \varphi(s) = \sum_{s \in S} \sum_{k=0}^{\infty} \left[ P^{\pi}(s | s_0, k) \sum_{a \in A} (Q^{\pi}(s, a) \nabla_{\theta} \pi_{\theta}(a | s)) \right]\end{aligned}$$

## Policy gradient theorem: the proof

- ▶ if we define  $d^\pi(s) := \sum_{k=0}^{\infty} P^\pi(s | s_0, k)$  as the non-normalised visitation probabilities of state  $s$  (starting from  $s_0$ ), the following

$$\nabla_\theta V^\pi(s_0) = \sum_{s \in \mathcal{S}} \sum_{k=0}^{\infty} \left[ P^\pi(s | s_0, k) \sum_{a \in \mathcal{A}} (Q^\pi(s, a) \nabla_\theta \pi_\theta(a | s)) \right]$$

becomes

$$\nabla_\theta V^\pi(s_0) = \sum_{s \in \mathcal{S}} d^\pi(s) \sum_{a \in \mathcal{A}} (Q^\pi(s, a) \nabla_\theta \pi_\theta(a | s))$$

- ▶ in the episodic case,  $d^\pi(s)$  the average length of an episode; in the continuing case it is 1.
- ▶ a Markov chain is ergodic  $\rightarrow d^\pi(s)$  a unique (non-normalised) stationary visiting prob. regardless of  $s_0$


# REINFORCE

- ▶ in order to make an unbiased estimation, the gradient can be further written as<sup>11</sup>:

$$\begin{aligned}\nabla_{\theta} V^{\pi}(s_0) &= \sum_{s \in S} d^{\pi}(s) \sum_{a \in A} (Q^{\pi}(s, a) \nabla_{\theta} \pi_{\theta}(a | s)) \\ &\propto \sum_{s \in S} \mu(s) \sum_{a \in A} \left( \pi_{\theta}(a | s) Q^{\pi}(s, a) \frac{\nabla_{\theta} \pi_{\theta}(a | s)}{\pi_{\theta}(a | s)} \right) \\ &= \mathbb{E}_{s \sim \mu, a \sim \pi} [Q^{\pi}(s, a) \nabla_{\theta} \ln \pi_{\theta}(a | s)] \\ &= \mathbb{E}_{s \sim \mu, a \sim \pi} [G_t \nabla_{\theta} \ln \pi_{\theta}(a | s)] \\ &\text{where } \mathbb{E}_{s \sim d^{\pi}, a \sim \pi} [G_t | s, a] = Q^{\pi}(s, a)\end{aligned}$$

- ▶ so the gradient update is  $\theta_{t+1} = \theta_t + \alpha G_t \nabla_{\theta} \ln \pi_{\theta}(a | s)$
- ▶ it is Monte Carlo Policy Gradient as REINFORCE uses the complete return from time  $t$ , which includes all future rewards up until the end for the episode

---

<sup>11</sup>Ronald J Williams. "Simple statistical gradient-following algorithms for connectionist reinforcement learning". In: *Machine learning* 8.3-4 (1992).  40/73



---

**Algorithm 4** Policy gradient with Monte-Carlo simulator

---

- 1: Initialize  $\theta$
  - 2: **repeat**
  - 3:     Sample trajectories  $\{\tau_i\}$  with horizon  $H$  using  $\pi_\theta(a|s)$
  - 4:      $G_{i,t} \leftarrow \sum_{t=t'}^H \gamma^{t-t'} R(s_{i,t}, a_{i,t})$
  - 5:      $V_t \leftarrow \frac{1}{M} \sum_{i=1}^M G_{i,t}$
  - 6:      $A(s_{i,t}, a_{i,t}) \leftarrow G_{i,t} - V_t$
  - 7:      $\Delta \leftarrow \sum_{i,t} \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) A(s_{i,t}, a_{i,t})$
  - 8:      $\theta \leftarrow \theta + \alpha \Delta$
  - 9: **until** convergence
- 

typically replace  $Q^{\pi_\theta}(s, a)$  with advantage function

$$A^{\pi_\theta}(s, a) := Q^{\pi_\theta}(s, a) - V^{\pi_\theta}(s).$$

in order to reduce the variance of the gradient

# Table of content

- ▶ Recap
  - ▶ MDPs
  - ▶ Value iterations and policy iterations
  - ▶ Tabular Q-Learning
- ▶ Policy approaches
  - ▶ Markov chains
  - ▶ Policy gradient
- ▶ **Computational learning theory**
  - ▶ PAC learning concepts
  - ▶ Learning bound for finite  $H$
- ▶ Theoretical analysis
  - ▶ Approximate dynamic programming
  - ▶ Performance bounds
  - ▶ Sample complexity

# Probably Approximately Correct (PAC) learning

- ▶ *learnability* is a key concept in ML:
  - ▶ what concepts can be learned?
  - ▶ how efficient is a particular learning method?
  - ▶ what is inherently hard to learn?
  - ▶ how many examples are needed in order to learn a concept successfully?
  - ▶ is there any generic model and theory about learnability?
- ▶ the Probably Approximately Correct (PAC) learning framework [Val84] is designed to answer the following two critical questions:
  1. sample complexity — how many training examples do we need to converge to a successful hypothesis with a high probability?
  2. computational complexity — how much computational effort is needed to converge to a successful hypothesis with a high probability?

## PAC: definition and notation

- ▶  $X$  as the set of all possible instances or examples, e.g., the set of images containing faces or non-faces classes
- ▶  $Y$  is the concept class, a set of target concepts  $y$ , e.g., face, non-face
- ▶ consider  $f : X \rightarrow Y = \{0,1\}$  the target concept to learn
- ▶ define  $D$  the target distribution, a fixed probability distribution over  $X$ . We shall make sure that the training and test examples are drawn according to  $D$
- ▶ define a set of training samples as  $S$  and a set of concept hypotheses  $H$ , e.g., the set of all linear classifiers

the learning problem is to, given a limited set of sample  $S$ , learn a hypothesis  $h : X \rightarrow Y \in H$  that approximating  $f$ . note that  $f$  may be in  $H$  or may not.

## PAC: definition and notation

- ▶ we then define two types of errors in order to understand the approximation
- ▶ true error or generalisation error of  $h$  is given as
$$R(h) = \Pr_{x \sim D} [h(x) \neq f(x)] = E_{x \sim D} [\mathbf{1}_{h(x) \neq f(x)}],$$
- ▶ whereas the average error of  $h$  on the training sample  $S$  is given according to the empirical distribution  $\hat{D}$  for the set  $S$ :
$$\hat{R}_S(h) = \Pr_{x \sim \hat{D}} [h(x) \neq f(x)] = E_{x \sim \hat{D}} [\mathbf{1}_{h(x) \neq f(x)}] = \frac{1}{m} \sum_{i=1}^m \mathbf{1}_{h(x_i) \neq f(x_i)}.$$
- ▶ note that  $R(h) = E_{S \sim D^m} [\hat{R}_S(h)]$ , where  $D^m$  is a distribution of sampling  $S$  according to  $D$ .

**Definition PAC Learning** [Val84]: A concept class  $Y$  is PAC-learnable if there exists an algorithm  $L(S) \rightarrow h$  such that:

- ▶ for all  $y \in Y$  and  $\delta > 0$  and  $\varepsilon > 0$  and all distributions  $D$ ,

$$\Pr_{S \sim D^m} [R(h) \leq \varepsilon] \geq 1 - \delta$$

- ▶ for samples  $S$  of size  $m \geq \text{poly}(1/\varepsilon, 1/\delta)$ , where  $\text{poly}()$  is a polynomial function.

PAC makes use of  $\delta > 0$  to define the confidence  $1 - \delta$  (*probabilistically*) and  $\varepsilon > 0$  the accuracy  $1 - \varepsilon$  (*approximate correct*)

A concept class  $Y$  is thus PAC-learnable if the hypothesis returned by the algorithm after observing a number of points polynomial in  $1/\varepsilon$  and  $1/\delta$  is approximately correct (error at most  $\varepsilon$ ) with high probability (at least  $1 - \delta$ )

## Learning bound for finite $H$ - consistent case

- ▶ **theorem:** let  $H$  be a finite set of functions from  $X$  to  $\{0, 1\}$  and  $L$  an algorithm that for any target concept  $y \in Y$  and sample  $S$  returns a *consistent*<sup>12</sup> hypothesis  $h_S : R_S(h_S) = 0$ . then, for any  $\delta > 0$ , with probability at least  $1 - \delta$


$$R(h_S) \leq \underbrace{\frac{1}{m}(\log |H| + \log \frac{1}{\delta})}_{\epsilon} \quad \text{generalisation bound}$$

the upper bound increases with  $\log |H|$  or the related term  $\log_2 |H|$ , which can be interpreted as the number of bits needed to represent  $H$

- ▶ Equivalently  $P_S D^m [R(h_S) \leq \epsilon] \geq 1 - \delta$  holds if

$$m \geq \frac{1}{\epsilon} (\log |H| + \log \frac{1}{\delta})$$

---

<sup>12</sup>a hypothesis set is consistent if it admits no error on training sample  $S$  

## Learning bound for finite $H$ - consistent case

- **proof:** for any  $\epsilon > 0$ , we define  $H_\epsilon = \{h \in H : R(h) > \epsilon\}$   
We then have:

$$\begin{aligned} & Pr[\exists h \in H_\epsilon : \hat{R}_S(h) = 0] && \Leftrightarrow R > \epsilon \cup R_S(h) = 0 \\ = & Pr[\hat{R}_S(h_1) = 0 \cup \dots \cup \hat{R}_S(h_{|H_\epsilon|}) = 0] \\ \leq & \sum_{h \in H_\epsilon} Pr[R_S(h) = 0] && \Leftrightarrow \text{union bound} \\ \leq & \sum_{h \in H_\epsilon} (1 - \epsilon)^m && \Leftrightarrow \text{no error in any } m \text{ samples} \\ \leq & |H|(1 - \epsilon)^m && \Leftrightarrow |H_\epsilon| \leq |H| \\ \leq & |H|e^{-m\epsilon} && \Leftrightarrow 1 - \epsilon \leq e^{-\epsilon} \end{aligned}$$



## Learning bound for finite $H$ : inconsistent Case

- ▶ no  $h \in H$  is a consistent hypothesis
- ▶ the typical case in practice: difficult problems, complex concept class
- ▶ but, inconsistent hypotheses with a small number of errors on the training set can be useful
- ▶ need a more powerful tool: Hoeffding's inequality

# Hoeffding's inequality

- ▶ **corollary:** for any  $\epsilon > 0$  and any hypothesis  $h : X \rightarrow \{0, 1\}$  the following inequality holds:

$$\Pr[|R(h) - \hat{R}(h)| \geq \epsilon] \leq 2e^{-2m\epsilon^2}$$

- ▶ this is due to:

$$\Pr[R(h) - \hat{R}(h) \geq \epsilon] \leq e^{-2m\epsilon^2}$$

$$\Pr[\hat{R}(h) - R(h) \geq \epsilon] \leq e^{-2m\epsilon^2}$$

- ▶ proof can be derived directly from Hoeffding's inequality [MRT18]

## Learning bound for finite $H$ : inconsistent Case

- ▶ **theorem:**  $H$  is a finite hypothesis set. for any  $\delta > 0$ , with probability at least  $1 - \delta$ ,

$$\forall h \in H, R(h) \leq \hat{R}_S(h) + \sqrt{\frac{\log |H| + \log \frac{2}{\delta}}{2m}}$$

- ▶ **proof:** By the union bound, we have

$$\begin{aligned} & Pr[\max_{h \in H} |R(h) - \hat{R}_S(h)| \geq \epsilon] \\ &= Pr[|R(h_1) - \hat{R}_S(h_1)| \geq \epsilon \cup \dots \cup |R(h_{|H|}) - \hat{R}_S(h_{|H|})| \geq \epsilon] \\ &\leq \sum_{h \in H} Pr[|R(h) - \hat{R}_S(h)| \geq \epsilon] \quad \Leftarrow \text{union bound} \\ &\leq 2|H|e^{-2m\epsilon^2} \quad \Leftarrow \text{apply previous corollary} \end{aligned}$$

## Estimation and approximation errors

- ▶  $H$  is a set of functions mapping from  $X$  to  $\{0, 1\}$ . The excess error (generalised error against Bayesian error<sup>13</sup>  $R^*$ ) of a  $h \in H$  can be decomposed as follows:

$$R(h_S) - R^* = \underbrace{(R(h_S) - \inf_{h \in H} R(h))}_{\text{estimation}} + \underbrace{(\inf_{h \in H} R(h) - R^*)}_{\text{approximation}}$$

- ▶ *estimation error*: it depends on the hypothesis  $h_S$  selected. it measures the error due to the samples
- ▶ *approximation error*: it measures how well the Bayes error can be approximated using  $H$

---

<sup>13</sup>The infimum of the errors achieved by any measurable functions:  
 $\min\{P[0|x], P[1|x]\}$

# Table of content

- ▶ Recap
  - ▶ MDPs
  - ▶ Value iterations and policy iterations
  - ▶ Tabular Q-Learning
- ▶ Policy approaches
  - ▶ Markov chains
  - ▶ Policy gradient
- ▶ Computational learning theory
  - ▶ PAC learning concepts
  - ▶ Learning bound for finite  $H$
- ▶ **Theoretical analysis**
  - ▶ Case study: approximate dynamic programming
  - ▶ Performance bounds - function approximation
  - ▶ Sample complexity - sampling-based

# Approximate methods

- ▶ when the state space is huge, approximate the value function by (parameterised) function approximators, e.g., linear model or neural networks
- ▶ In general, ADP (Approximate DP) approximates the optimal value function  $V^*$  by finding  $V$  in some function space  $\mathcal{V}$  as

$$V = \arg \min_{U \in \mathcal{V}} \text{distance}(V^*, U)$$

with some distance metric and use the approximator  $V$  to generate a greedy policy as

$$\pi(s) = \arg \max_a R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V(s).$$

- ▶ Most modern RL algorithms can be viewed as solving ADP by sampling the transition  $\Pr(s'|s, a)$  and rewards  $r(s, a, s')$
- ▶ Unlike DP, ADP algorithms may not converge as the composition of Bellman operator  $T$  and projection onto the function space  $\mathcal{V}$  is not a contraction general, resulting in oscillating and even divergent behavior like Q-learning
- ▶ Nevertheless, we can analyze its **performance bounds** and **sample complexity bounds**

## ADP performance bounds

- ▶ DP is guaranteed to converge and consequently leads to convergence of RL with tabular MDP
- ▶ Unlike DP, not all ADP algorithms converge and ADP may not find the actual (optimal) value functions for  $V^*$  is not necessarily in  $\mathcal{V}$
- ▶ The performance gap of a policy  $\pi$  generated from  $V$  and the actual optimal policy  $\pi^*$  has certain theoretic upper bounds
- ▶ One is given by difference between the optimal value function and the approximator  $V$  itself

$$\underbrace{\|V^* - V^\pi\|_\infty}_{\text{performance gap}} \leq \frac{2\gamma}{1-\gamma} \underbrace{\|V^* - V\|_\infty}_{\text{approximation error}}.$$

This follows from

$$\begin{aligned} \|V^* - V^\pi\|_\infty &\leq \|V^* - T^\pi V\|_\infty + \|T^\pi V - T^\pi V^\pi\|_\infty \\ &\leq \|TV^* - TV\|_\infty + \gamma \|V - V^\pi\|_\infty \\ &\leq \gamma \|V^* - V\|_\infty + \gamma (\|V - V^*\|_\infty + \|V^* - V^\pi\|_\infty) \\ &\leq \frac{2\gamma}{1-\gamma} \|V^* - V\|_\infty. \end{aligned} \tag{12}$$

# AVI and API

- ▶ this approximation error defined bound motivates algorithms such as AVI and API to minimise the approximation error
- ▶ formally, the approximate value iteration (AVI) is written as

$$\text{(AVI)} \quad V_{k+1} = \text{Proj}_{\mathcal{V}} TV_k$$

with projection  $\text{Proj}_{\mathcal{V}} V = \arg \min_{V' \in \mathcal{V}} \|V - V'\|$  for certain norm  $\|\cdot\|$ ; and

- ▶ the approximate policy iteration (API) is written as

$$\text{(API)} \quad \pi_{k+1}(s) = \arg \max_a R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) V(s'),$$

with  $V \approx V^{\pi_k}$ .

- ▶ in particular, the value function  $V$  in API is obtained through approximation in a function space  $\mathcal{V}$
- ▶ expect that if  $V$  is close to  $V^*$  then the policy  $\pi$  will be close to optimal



## Bounds by Bellman residual

- Apart from bound by approximation error, it is also possible to bound the performance with the Bellman residual as [WB93]

$$\| \underbrace{V^* - V^\pi}_{\text{performance gap}} \|_\infty \leq \frac{2}{1-\gamma} \| \underbrace{TV - V}_{\text{Bellman residual}} \|_\infty, \quad (13)$$

where  $T$  is the Bellman optimality operator. This follows by combing

$$\begin{aligned} \|V^* - V\|_\infty &\leq \|V^* - TV\|_\infty + \|TV - V\|_\infty \\ &\leq \gamma \|V^* - V\|_\infty + \|TV - V\|_\infty \leq \frac{1}{1-\gamma} \|TV - V\|_\infty \end{aligned}$$

and

$$\begin{aligned} \|V - V^\pi\|_\infty &\leq \|V - TV\|_\infty + \|TV - V^\pi\|_\infty \\ &\leq \|TV - V\|_\infty + \gamma \|V - V^\pi\|_\infty \quad (\text{By } TV = T^\pi V) \\ &\leq \frac{1}{1-\gamma} \|TV - V\|_\infty \end{aligned}$$

# Bellman residual minimisation

- ▶ this Bellman residual bound motivates one to minimise the Bellman residual, which leads to Bellman residual minimisation

$$\text{(BRM)} \quad \min_{V \in \mathcal{V}} \|TV - V\|$$

for some norm  $\|\cdot\|$

## Performance bounds of AVI

- ▶ the role of  $V_k$  in AVI is similar to that of the **target network** in DQN:
- ▶ AVI bound [BT96]: after  $K$  iterations, we have

$$\|V^* - V^{\pi_K}\|_\infty \leq \frac{2\gamma}{(1-\gamma)^2} \max_{0 \leq k \leq K} \|TV_k - V_{k+1}\|_\infty + \frac{2\gamma^{K+1}}{1-\gamma} \|V^* - V_0\|_\infty.$$

In particular, if

$$\tilde{V} = \text{Proj}_{\mathcal{V}} T\tilde{V}$$

with  $\tilde{\pi}$  being a greedy policy with respect to  $R + \gamma P\tilde{V}$ , then

$$\|V^* - V^{\tilde{\pi}}\|_\infty \leq \frac{2}{(1-\gamma)^2} \inf_{V \in \mathcal{V}} \|V^* - V\|_\infty.$$

## Performance bounds of AVI: the proof

- ▶ by letting  $\varepsilon = \max_{0 \leq k \leq K} \|TV_k - V_{k+1}\|_\infty$ , we can derive

$$\begin{aligned}\|V^* - V_{k+1}\|_\infty &\leq \|TV^* - TV_k\| + \|TV_k - V_{k+1}\|_\infty \\ &\leq \gamma \|V^* - V_k\|_\infty + \varepsilon,\end{aligned}$$

and thus,

$$\begin{aligned}\|V^* - V_k\|_\infty &\leq (1 + \gamma + \dots + \gamma^{k-1})\varepsilon + \gamma^k \|V^* - V_0\|_\infty \\ &\leq \frac{1}{1 - \gamma} \varepsilon + \gamma^k \|V^* - V_0\|_\infty.\end{aligned}$$

The first result follows by combining (12)

- ▶ let the projection use the infinity norm, then the AVI is contractive with fixed point  $\tilde{V} = \text{Proj}_{\mathcal{Y}} T\tilde{V}$  and we can obtain

$$\|V^* - V\|_\infty \leq \|V^* - \text{Proj}_{\mathcal{Y}} V^*\|_\infty + \|\text{Proj}_{\mathcal{Y}} V^* - \tilde{V}\|_\infty$$

with  $\|\text{Proj}_{\mathcal{Y}} V^* - \tilde{V}\|_\infty = \|\text{Proj}_{\mathcal{Y}} TV^* - \text{Proj}_{\mathcal{Y}} T\tilde{V}\| \leq \gamma \|V^* - \tilde{V}\|_\infty$ .  
the second result then follows from using (12).

# Performance bounds of API

- ▶ API bound [BT96]: the asymptotic performance bound is

$$\limsup_{k \rightarrow \infty} \|V^* - V^{\pi_k}\|_{\infty} \leq \frac{2\gamma}{(1-\gamma)^2} \limsup_{k \rightarrow \infty} \|V_k - V^{\pi_k}\|_{\infty}.$$

## Performance bounds of API: the proof

- ▶ let  $e_k = V_k - V^{\pi_k}$  denote the approximation error,  $g_k = V^{\pi_{k+1}} - V^{\pi_k}$  the performance gain and  $l_k = V^* - V^{\pi_k}$  the loss of using  $\pi_k$  instead of  $\pi^*$
- ▶ we can show that the next policy cannot be much worse than the current one as

$$g_k \geq -\gamma(I - \gamma P^{\pi_{k+1}})^{-1}(P^{\pi_{k+1}} - P^{\pi_k})e_k.$$

- ▶ the loss at the next iteration is bounded by the current loss as

$$l_{k+1} \leq \gamma P^{\pi^*} l_k + f_k$$

where  $f_k = \gamma[P^{\pi_{k+1}}(I - \gamma P^{\pi_{k+1}})^{-1}(I - \gamma P^{\pi_k}) - P^{\pi^*}]e_k$

- ▶ by taking the limit on both sides, we can obtain

$$\begin{aligned}(I - \gamma P^{\pi^*}) \limsup_{k \rightarrow \infty} l_k &\leq \limsup_{k \rightarrow \infty} f_k \\ \limsup_{k \rightarrow \infty} l_k &\leq (I - \gamma P^{\pi^*})^{-1} \limsup_{k \rightarrow \infty} f_k.\end{aligned}$$

This leads to

$$\begin{aligned}\limsup_{k \rightarrow \infty} \|l_k\|_\infty &\leq \frac{\gamma}{1 - \gamma} \|P^{\pi_{k+1}}(I - \gamma P^{\pi_{k+1}})^{-1}(I - \gamma P^{\pi_k}) - P^{\pi^*}\|_\infty \|e_k\|_\infty \\ &\leq \frac{\gamma}{1 - \gamma} \left(\frac{1 + \gamma}{1 - \gamma} + 1\right) \|e_k\|_\infty + \frac{2\gamma}{(1 - \gamma)^2} \|e_k\|_\infty,\end{aligned}$$

which validates the result

## Performance bounds of BRM

- BRM [WB93]: if  $V_{\text{BRM}} = \arg \min_{V \in \mathcal{V}} \|TV - V\|_{\infty}$ ,

$$\|V^* - V_{\text{BRM}}^{\pi}\| \leq \frac{2(1 + \gamma)}{1 - \gamma} \inf_{V \in \mathcal{V}} \|V^* - V\|_{\infty}.$$

**Proof.** Note that

$$\begin{aligned} \|TV - V\|_{\infty} &\leq \|TV - TV^*\|_{\infty} + \|V^* - V\|_{\infty} \\ &\leq (1 + \gamma) \|V^* - V\|_{\infty}. \end{aligned}$$

The bound follows by combining

$$\begin{aligned} \|TV_{\text{BRM}}^{\pi} - V_{\text{BRM}}^{\pi}\|_{\infty} &= \inf_{V \in \mathcal{V}} \|TV - V\|_{\infty} \\ &\leq (1 + \gamma) \inf_{V \in \mathcal{V}} \|V^* - V\|_{\infty} \end{aligned}$$

and (13).

## Sample-based ADP: sample complexity

- ▶ RL essentially solve ADP with sampled transition and rewards
- ▶ from statistical learning, the *prediction error* of RL comes from two sources:
  1. **approximation error**, error due to the projection operation;
  2. **estimation error**, since both the Bellman operator and projection are evaluated with samples.
- ▶ error propagation of sample-based ADP can be analyzed by using tools such as McDiarmid's inequality, if

$$\sup_{x_1, \dots, x_n, x'_i} |f(x_1, \dots, x_i, \dots, x_n) - f(x_1, \dots, x'_i, \dots, x_n)| \leq c_i,$$

we have

$$\Pr(|f(x_1, \dots, x_n) - \mathbb{E}[f(x_1, \dots, x_n)]| \geq \varepsilon) \leq \exp\left(-\frac{2\varepsilon}{\sum_{i=1}^n c_i}\right).$$

- ▶ combing the sample error and performance bounds before, the generalisation bound (performance guarantee) of sample-based ADPs can be derived



## Sample complexity: sampling-based AVI

- ▶ consider following setup [MS08]: sample i.i.d.  $n$  states  $s^{(i)} \sim \mu$ , and from each state-action pair  $s^{(i)}, a$ , generate  $m$  one-step transition samples from a simulator  $s_a^{(i,j)} \sim \Pr(\cdot | s^{(i)}, a)$
- ▶ iterate AVI with **fitted value functions**  $K$  times:

$$V_{k+1} = \arg \min_{V \in \mathcal{V}} \sum_{i=1}^n |V(s^{(i)}) - \max_a [r(s^{(i)}, a) + \frac{\gamma}{m} \sum_{j=1}^m V_k(s_a^{(i,j)})]|^2.$$

## Sample complexity: sampling-based AVI

- ▶ With probability at least  $1 - \delta$ ,

$$\begin{aligned} \|V^* - V^{\pi_K}\|_\infty &\leq \frac{2\gamma}{(1-\gamma)^2} C^{1/p} d(T\mathcal{V}, \mathcal{V}) + O(\gamma^k) \\ &\quad + O\left(\frac{V(\mathcal{V}) \log(1/\delta)}{n}\right)^{1/4} + O\left(\frac{\log(1/\delta)}{m}\right)^{1/2}, \end{aligned}$$


where

$$d(T\mathcal{V}, \mathcal{V}) := \sup_{V \in \mathcal{V}} \inf_{V' \in \mathcal{V}} \|TV - V'\|_{2,\mu}, \text{ with } \|V\|_{2,\mu} = \left(\sum_x \mu(x) V(x)^2\right)^{1/2}$$

measures the Bellman residual of the space  $\mathcal{V}$ , the constant  $C$  satisfies  $1 \leq C \leq \Pr(\cdot|s, a)/\mu(\cdot)$  for any  $s$  and  $a$ , and  $V(\mathcal{V})$  is the capacity measure of  $\mathcal{V}$  (i.e., pseudo-dimension<sup>14</sup>).

- ▶ In addition to the **AVI** with fitted value iteration [MS08], PAC bounds for finite-time analysis has also been developed for **API** such as LSPTD/LSPI by [LGM12] and BRM-based PI by [Mai+10]

---

<sup>14</sup>The Pseudo-dimension, also referred to the Pollard dimension, is a generalization of the VC-dimension to real-valued functions. 

## Deep Q learning and its sample complexity

- ▶ finite time bound of Q-learning with non-linear multi-layer ReLU unit and i.i.d. samples have been studied [YXW19]
- ▶ the major result states that, with high probability,

$$\|Q^{\pi_K} - Q^*\|_{1,\mu} \leq C \frac{\phi_{\mu,\sigma} \cdot \gamma}{(1-\gamma)^2} |\mathbb{A}| (\log n)^{1+2\xi^*} n^{(\alpha^*-1)/2} + \frac{4\gamma^{K+1}}{1-\gamma} \cdot \max_{s,a} R(s,a),$$

where  $n$  is the number of samples,  $C > 0$ ,  $\xi^*$  and  $\alpha^*$  is some constant,  $\phi_{\mu,\sigma}$  is related to concentration coefficients of underlying Markov chain

- ▶ this was further extended to non i.i.d. samples by [XG19]. With high probability, the bound decreases at rate  $1/\sqrt{K}$  with sufficiently network width  $m$  as

$$\frac{1}{K} \sum_{k=0}^K \mathbb{E}[(Q(s,a;\theta_k) - Q^*(s,a))^2] \leq O\left(\frac{1}{m^{1/6}} + \frac{1}{\sqrt{K}}\right)$$

## Concluding remarks

# References I



Richard Bellman. *Dynamic Programming*. Princeton University Press, 1957.



Dimitri P Bertsekas and John N Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, 1996.



Ronald A Howard. *Dynamic Programming and Markov Processes*. MIT Press, 1960.



Alessandro Lazaric, Mohammad Ghavamzadeh, and Rémi Munos. “Finite-sample analysis of least-squares policy iteration”. In: *The Journal of Machine Learning Research* 13 (2012), pp. 3041–3074.



Odalric-Ambrym Maillard et al. “Finite-sample analysis of Bellman residual minimization”. In: *Asian Conference on Machine Learning (ACML)*. 2010, pp. 299–314.

## References II



Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.



Rémi Munos and Csaba Szepesvári. “Finite-time bounds for fitted value iteration”. In: *Journal of Machine Learning Research* 9 (2008), pp. 815–857.



Rémi Munos. *Introduction to Reinforcement Learning and multi-armed bandits*. NETADIS Summer School, 2013.



Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.



Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

## References III



Richard Sutton and Andrew Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.



Richard S Sutton et al. “Policy gradient methods for reinforcement learning with function approximation”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 1999, pp. 1057–1063.



Leslie G Valiant. “A theory of the learnable”. In: *Communications of the ACM* 27.11 (1984), pp. 1134–1142.



Christopher John Cornish Hellaby Watkins. “Learning From Delayed Rewards”. PhD Thesis. University of Cambridge, 1989.

## References IV



Ronald J. Williams and Leemon C. Baird III. *Tight performance bounds on greedy policies based on imperfect value functions*. Tech. rep. NU-CCS-93-14, College of Computer Science, Northeastern University. 1993.



Ronald J Williams. “Simple statistical gradient-following algorithms for connectionist reinforcement learning”. In: *Machine learning* 8.3-4 (1992).



Shuang Wu and Jun Wang. *Decision making and AI: a white paper*. 2020.



Pan Xu and Quanquan Gu. “A finite-time analysis of q-learning with neural network function approximation”. In: *arXiv preprint arXiv:1912.04511* (2019).



## References V



Zhuoran Yang, Yuchen Xie, and Zhaoran Wang. “A theoretical analysis of deep Q-learning”. In: *arXiv preprint arXiv:1901.00137* (2019).