

# Lecture 4: Model Free Control <sup>2</sup>

Emma Brunskill

CS234 Reinforcement Learning.

Winter 2018

---

<sup>2</sup>Structure closely follows much of David Silver's Lecture 5. For additional reading please see SB Sections 5.2-5.4, 6.4, 6.5, 6.7

# Table of Contents

- 1 Generalized Policy Iteration
- 2 Importance of Exploration
- 3 Monte Carlo Control
- 4 Temporal Difference Methods for Control
- 5 Maximization Bias

# Class Structure

- Last time: Policy evaluation with no knowledge of how the world works (MDP model not given)
- This time: Control (making decisions) without a model of how the world works
- Next time: Value function approximation and Deep Q-learning

# Evaluation to Control

- Last time: how good is a specific policy?
  - Given no access to the decision process model parameters
  - Instead have to estimate from data / experience
- Today: how can we learn a good policy?

# Recall: Reinforcement Learning Involves

- Optimization
- Delayed consequences
- Exploration
- Generalization

# Today: Learning to Control Involves

- Optimization: Goal is to identify a policy with high expected rewards (similar to Lecture 2 on computing an optimal policy given decision process models)
- Delayed consequences: May take many time steps to evaluate whether an earlier decision was good or not
- Exploration: Necessary to try different actions to learn what actions can lead to high rewards

# Today: Model-free Control

- Generalized policy improvement
- Importance of exploration
- Monte Carlo control
- Model-free control with temporal difference (SARSA, Q-learning)
- Maximization bias

# Model-free Control Examples

- Many applications can be modeled as a MDP: Backgammon, Go, Robot locomotion, Helicopter flight, Robocup soccer, Autonomous driving, Customer ad selection, Invasive species management, Patient treatment
- For many of these and other problems either:
  - MDP model is unknown but can be sampled
  - MDP model is known but it is computationally infeasible to use directly, except through sampling



# On and Off-Policy Learning

- On-policy learning
  - Direct experience
  - Learn to estimate and evaluate a policy from experience obtained from following that policy
- Off-policy learning
  - Learn to estimate and evaluate a policy using experience gathered from following a different policy

# Table of Contents

- 1 Generalized Policy Iteration
- 2 Importance of Exploration
- 3 Monte Carlo Control
- 4 Temporal Difference Methods for Control
- 5 Maximization Bias

# Recall Policy Iteration

- Initialize policy  $\pi$
- Repeat:
  - Policy evaluation: compute  $V^\pi$
  - Policy improvement: update  $\pi$

$$\pi'(s) = \arg \max_a R(s, a) + \gamma \sum_{s' \in S} P(s'|s, a) V^\pi(s') = \arg \max_a Q^\pi(s, a) \quad (1)$$

- Now want to do the above two steps without access to the true dynamics and reward models
- Last lecture introduced methods for model-free policy evaluation

# Model-free Generalized Policy Improvement

- Given an estimate  $Q^{\pi_i}(s, a) \forall s, a$
- Update new policy

$$\pi_{i+1}(s) = \arg \max_a Q^{\pi_i}(s, a) \quad (2)$$

# Model-free Policy Iteration

- Initialize policy  $\pi$
- Repeat:
  - Policy evaluation: compute  $Q^\pi$
  - Policy improvement: update  $\pi$  given  $Q^\pi$
- May need to modify policy evaluation:
  - If  $\pi$  is deterministic, can't compute  $Q(s, a)$  for any  $a \neq \pi(s)$
- How to interleave policy evaluation and improvement?
  - Policy improvement is now using an estimated  $Q$

model-free  
← policy eval

# Table of Contents

- 1 Generalized Policy Iteration
- 2 Importance of Exploration**
- 3 Monte Carlo Control
- 4 Temporal Difference Methods for Control
- 5 Maximization Bias

# Policy Evaluation with Exploration

- Want to compute a model-free estimate of  $Q^\pi$
- In general seems subtle
  - Need to try all  $(s, a)$  pairs but then follow  $\pi$
  - Want to ensure resulting estimate  $Q^\pi$  is good enough so that policy improvement is a monotonic operator
- For certain classes of policies can ensure all  $(s, a)$  pairs are tried such that asymptotically  $Q^\pi$  converges to the true value

# $\epsilon$ -greedy Policies

- Simple idea to balance exploration and exploitation
- Let  $|A|$  be the number of actions
- Then an  $\epsilon$ -greedy policy w.r.t. a state-action value  $Q^\pi(s, a)$  is

$$\pi(a|s) = \begin{cases} a & \text{w. prob } \epsilon / |A| \\ \operatorname{argmax}_a Q^\pi(s, a) & \text{w. prob } 1 - \epsilon \end{cases}$$



# Monotonic<sup>19</sup> $\epsilon$ -greedy Policy Improvement

## Theorem

For any  $\epsilon$ -greedy policy  $\pi_i$ , the  $\epsilon$ -greedy policy w.r.t.  $Q^{\pi_i}$ ,  $\pi_{i+1}$  is a monotonic improvement  $V^{\pi_{i+1}} \geq V^{\pi_i}$

$$\begin{aligned} Q^{\pi}(s, \pi_{i+1}(s)) &= \sum_{a \in A} \pi_{i+1}(a|s) Q^{\pi_i}(s, a) \quad \text{randomly} \\ &= (\epsilon/|A|) \sum_{a \in A} Q^{\pi_i}(s, a) + (1-\epsilon) \max_a Q^{\pi_i}(s, a) \quad \text{greedy} \\ &= \epsilon/|A| \sum_a Q^{\pi_i}(s, a) + (1-\epsilon) \max_a Q^{\pi_i}(s, a) \quad \text{for any state} \\ &= \epsilon/|A| \sum_a Q^{\pi_i}(s, a) + (1-\epsilon) \frac{1-\epsilon}{1-\epsilon} \sum_a (\pi_i(a|s) - \epsilon/|A|) Q^{\pi_i}(s, a) \\ &\geq \epsilon/|A| \sum_a Q^{\pi_i}(s, a) + (1-\epsilon) \frac{1}{1-\epsilon} \sum_a (\pi_i(a|s) - \epsilon/|A|) Q^{\pi_i}(s, a) \\ &= \sum_a \pi_i(a|s) Q^{\pi_i}(s, a) = V^{\pi_i}(s) \end{aligned}$$

- Therefore  $V^{\pi_{i+1}} \geq V^{\pi_i}$  (from the policy improvement theorem)

<sup>19</sup>The theorem assumes that  $Q^{\pi_i}$  has been computed exactly.


# Monotonic<sup>21</sup> $\epsilon$ -greedy Policy Improvement

## Theorem

For any  $\epsilon$ -greedy policy  $\pi_i$ , the  $\epsilon$ -greedy policy w.r.t.  $Q^{\pi_i}$ ,  $\pi_{i+1}$  is a monotonic improvement  $V^{\pi_{i+1}} \geq V^{\pi}$

$$\begin{aligned} Q^{\pi}(s, \pi_{i+1}(s)) &= \sum_{a \in A} \pi_{i+1}(a|s) Q^{\pi_i}(s, a) \\ &= (\epsilon/|A|) \sum_{a \in A} Q^{\pi_i}(s, a) + (1 - \epsilon) \max_a Q^{\pi_i}(s, a) \\ &= (\epsilon/|A|) \sum_{a \in A} Q^{\pi_i}(s, a) + (1 - \epsilon) \max_a Q^{\pi_i}(s, a) \frac{1 - \epsilon}{1 - \epsilon} \\ &= (\epsilon/|A|) \sum_a Q^{\pi_i}(s, a) + (1 - \epsilon) \max_a Q^{\pi_i}(s, a) \sum_a \frac{\pi_i(a|s) - \frac{\epsilon}{|A|}}{1 - \epsilon} \\ &\geq \frac{\epsilon}{|A|} \sum_{a \in A} Q^{\pi_i}(s, a) + (1 - \epsilon) \sum_a \frac{\pi_i(a|s) - \frac{\epsilon}{|A|}}{1 - \epsilon} Q^{\pi_i}(s, a) \\ &= \sum_a \pi_i(a|s) Q^{\pi_i}(s, a) = V^{\pi_i}(s) \end{aligned}$$

- Therefore  $V^{\pi_{i+1}} \geq V^{\pi}$  (from the policy improvement theorem)

<sup>21</sup>The theorem assumes that  $Q^{\pi_i}$  has been computed exactly. 

# Greedy in the Limit of Infinite Exploration (GLIE)

## Definition of GLIE

- All state-action pairs are visited an infinite number of times

$$\lim_{i \rightarrow \infty} N_i(s, a) \rightarrow \infty$$

- Behavior policy converges to greedy policy

$$\lim_{i \rightarrow \infty} \pi(a|s) \Rightarrow \arg \max_a q(s, a) \quad \text{prob} \rightarrow$$

↙ optimal state-action v

- A simple GLIE strategy is  $\epsilon$ -greedy where  $\epsilon$  is reduced to 0 with the following rate:  $\epsilon_i = 1/i$

# Table of Contents

- 1 Generalized Policy Iteration
- 2 Importance of Exploration
- 3 Monte Carlo Control**
- 4 Temporal Difference Methods for Control
- 5 Maximization Bias

# Monte Carlo Online Control / On Policy Improvement

- 
- 1: Initialize  $Q(s, a) = 0$ ,  $Returns(s, a) = 0 \forall (s, a)$ , Set  $\epsilon = 1$ ,  $k = 1$
  - 2:  $\pi_k = \epsilon\text{-greedy}(Q)$  // Create initial  $\epsilon\text{-greedy}$  policy
  - 3: **loop**
  - 4:   Sample  $k$ -th episode  $(s_{k1}, a_{k1}, r_{k1}, s_{k2}, \dots, s_T)$  given  $\pi_k$
  - 5:   **for**  $t = 1, \dots, T$  **do**
  - 6:     **if** First visit to  $(s, a)$  in episode  $k$  **then**   *could use discount  $\gamma$  on rewards*
  - 7:       Append  $\sum_{j=t}^T r_{kj}$  to  $Returns(s_t, a_t)$
  - 8:        $Q(s_t, a_t) = \text{average}(Returns(s_t, a_t))$
  - 9:     **end if**
  - 10:   **end for**
  - 11:    $k = k + 1$ ,  $\epsilon = 1/k$     *$\epsilon\text{-greedy}(Q) \rightarrow$  compute an  $\epsilon\text{-greedy}$  policy*
  - 12:    $\pi_k = \epsilon\text{-greedy}(Q^\pi)$  // Policy improvement
  - 13: **end loop**
-

## Theorem

GLIE Monte-Carlo control converges to the optimal state-action value<sup>a</sup> function  $Q(s, a) \rightarrow q(s, a)$

---

<sup>a</sup> $v(s)$  and  $q(s, a)$  without any additional subscripts are used to indicate the optimal state and state-action value function, respectively.

# Model-free Policy Iteration

- Initialize policy  $\pi$
- Repeat:
  - Policy evaluation: compute  $Q^\pi$
  - Policy improvement: update  $\pi$  given  $Q^\pi$
- What about TD methods?

# Table of Contents

- 1 Generalized Policy Iteration
- 2 Importance of Exploration
- 3 Monte Carlo Control
- 4 Temporal Difference Methods for Control**
- 5 Maximization Bias



# Model-free Policy Iteration with TD Methods

- Use temporal difference methods for policy evaluation step
- Initialize policy  $\pi$
- Repeat:
  - Policy evaluation: compute  $Q^\pi$  using temporal difference updating with  $\epsilon$ -greedy policy
  - Policy improvement: Same as Monte carlo policy improvement, set  $\pi$  to  $\epsilon$ -greedy ( $Q^\pi$ )

TD learning

$$V^\pi(s) = V^\pi(s) + \alpha \left( \underbrace{r + \gamma V^\pi(s')}_{\text{bootstrapped}} - \underbrace{V^\pi(s)}_{\text{old estimate}} \right)$$

Handwritten annotations: "sampling" with arrows pointing to  $s$  and  $s'$ ; "learning ref" with an arrow pointing to  $V^\pi(s)$  on the left; "bootstrapped" under the bracketed term; "old estimate" under the bracketed term.

# General Form of SARSA Algorithm

on policy

- 1: Set initial  $\epsilon$ -greedy policy  $\pi$ ,  $t = 0$ , initial state  $s_t = s_0$   $Q(s, a) = 0$
- 2: Take  $a_t \sim \pi(s_t)$  // Sample action from policy  $\forall s, a$
- 3: Observe  $(r_t, s_{t+1})$
- 4: **loop**
- 5: Take action  $a_{t+1} \sim \pi(s_{t+1})$  action that would be taken
- 6: Observe  $(r_{t+1}, s_{t+2})$
- 7: Update Q given  $(s_t, a_t, r_t, s_{t+1}, a_{t+1})$ :  
$$Q^\pi(s_t, a_t) \leftarrow Q^\pi(s_t, a_t) + \alpha (r_t + \gamma Q^\pi(s_{t+1}, a_{t+1}) - Q^\pi(s_t, a_t))$$
 $\uparrow$   
learning rate
- 8: Perform policy improvement:  
 $\pi \leftarrow \epsilon\text{-greedy}(Q^\pi)$
- 9:  $t = t + 1$
- 10: **end loop**

- What are the benefits to improving the policy after each step?

# Convergence Properties of SARSA

## Theorem

Sarsa for finite-state and finite-action MDPs converges to the optimal action-value,  $Q(s, a) \rightarrow q(s, a)$ , under the following conditions:

- 1 The policy sequence  $\pi_t(a|s)$  satisfies the condition of GLIE
- 2 The step-sizes  $\alpha_t$  satisfy the Robbins-Munro sequence such that

$$\sum_{t=1}^{\infty} \alpha_t = \infty$$

$$\sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

# Recall: Off Policy, Policy Evaluation

- Given data from following a behavior policy  $\pi_b$  can we estimate the value  $V^{\pi_e}$  of an alternate policy  $\pi_e$ ?
- Neat idea: can we learn about other ways to do things different than what we actually did?
- Discussed how to do this for Monte Carlo evaluation
- Used Importance Sampling
- First see how to do off policy evaluation with TD

# Importance Sampling for Off Policy TD (Policy Evaluation)

- Recall the Temporal Difference (TD) algorithm which is used to incremental model-free evaluation of a policy  $\pi_b$ . Precisely, given a state  $s_t$ , an action  $a_t$  sampled from  $\pi_b(s_t)$  and the observed reward  $r_t$  and next state  $s_{t+1}$ , TD performs the following update:

$$V^{\pi_b}(s_t) = V^{\pi_b}(s_t) + \alpha(r_t + \gamma V^{\pi_b}(s_{t+1}) - V^{\pi_b}(s_t)) \quad (3)$$

- Now want to use data generated from following  $\pi_b$  to estimate the value of different policy  $\pi_e$ ,  $V^{\pi_e}$
- Change TD target  $r_t + \gamma V(s_{t+1})$  to weight target by single importance sample ratio
- New update:

*likelihood under eval  $\pi_e$*   
 $\frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)}$

$$V^{\pi_e}(s_t) = V^{\pi_e}(s_t) + \alpha \left[ \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)} (r_t + \gamma V^{\pi_e}(s_{t+1}) - V^{\pi_e}(s_t)) \right] \quad (4)$$

# Importance Sampling for Off Policy TD Cont.

- Off Policy TD Update:

$$\prod_{t=1}^T \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)} \quad \frac{1}{1/2}$$

$$V^{\pi_e}(s_t) = V^{\pi_e}(s_t) + \alpha \left[ \frac{\pi_e(a_t|s_t)}{\pi_b(a_t|s_t)} (r_t + \gamma V^{\pi_e}(s_{t+1}) - V^{\pi_e}(s_t)) \right] \quad (5)$$

- Significantly lower variance than MC IS. (Why?)

Yes, only 1 importance weight instead of product

- Does  $\pi_b$  need to be the same at each time step?

No, just need to know it

- What conditions on  $\pi_b$  and  $\pi_e$  are needed for off policy TD to converge to  $V^{\pi_e}$ ?  $\pi_b$  has the same support

$$\pi_b(a|s) > 0 \text{ whenever } \pi_e(a|s) \cdot V(s) > 0$$

# Q-Learning: Learning the Optimal State-Action Value

- Just saw how to use off policy TD to evaluate any particular policy  $\pi_e$
- Can we estimate the value of the optimal policy  $\pi^*$  without knowledge of what  $\pi^*$  is?
- Yes! Q-learning
- Does not require importance sampling
- Key idea: Maintain state-action  $Q$  estimates and use to bootstrap—use the value of the best future action
- Recall Sarsa

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha((r_t + \gamma Q(s_{t+1}, a_{t+1})) - Q(s_t, a_t)) \quad (6)$$

- Q-learning:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha((r_t + \gamma \max_{a'} Q(s_{t+1}, a')) - Q(s_t, a_t)) \quad (7)$$

# Off-Policy Control Using Q-learning

- In the prior slide assumed there was some  $\pi_b$  used to act
- $\pi_b$  determines the actual rewards received
- Now consider how to improve the behavior policy (policy improvement)
- Let behavior policy  $\pi_b$  be  $\epsilon$ -greedy with respect to (w.r.t.) current estimate of the optimal  $q(s, a)$



# Q-Learning with $\epsilon$ -greedy Exploration

1: Initialize  $Q(s, a), \forall s \in S, a \in A$   $t = 0$ , initial state  $s_t = s_0$

2: Set  $\pi_b$  to be  $\epsilon$ -greedy w.r.t.  $Q$

3: **loop**

4: Take  $a_t \sim \pi_b(s_t)$  // Sample action from policy

5: Observe  $(r_t, s_{t+1})$

6: Update  $Q$  given  $(s_t, a_t, r_t, s_{t+1})$ :

$$Q(s_t, a) = Q(s_t, a) + \alpha (r_t + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a))$$

7: Perform policy improvement: set  $\pi_b$  to be  $\epsilon$ -greedy w.r.t.  $Q$

8:  $t = t + 1$

9: **end loop**

- What conditions are sufficient to ensure that Q-learning with  $\epsilon$ -greedy exploration converges to optimal  $q$ ? *visit all  $s,a$  pairs  $\infty$  often*
- What conditions are sufficient to ensure that Q-learning with  $\epsilon$ -greedy exploration converges to optimal  $\pi^*$ ? *decay  $\epsilon$*

# Table of Contents

- 1 Generalized Policy Iteration
- 2 Importance of Exploration
- 3 Monte Carlo Control
- 4 Temporal Difference Methods for Control
- 5 Maximization Bias

# Maximization Bias<sup>39</sup>

- Consider single-state MDP ( $|S| = 1$ ) with 2 actions, and both actions have 0-mean random rewards, ( $\mathbb{E}(r|a = a_1) = \mathbb{E}(r|a = a_2) = 0$ ).

- Then  $Q(s, a_1) = Q(s, a_2) = 0 = V(s)$

- Assume there are prior samples of taking action  $a_1$  and  $a_2$

- Let  $\hat{Q}(s, a_1), \hat{Q}(s, a_2)$  be the finite sample estimate of  $Q$

- Assume using an unbiased estimator for  $Q$ : e.g.

$$\hat{Q}(s, a_1) = \frac{1}{n(s, a_1)} \sum_{i=1}^{n(s, a_1)} r_i(s, a_1)$$

ex.  
100 samples  
 $a_1$   
100 samples  
 $a_2$

- Let  $\hat{\pi} = \arg \max_a \hat{Q}(s, a)$  be the greedy policy w.r.t. the estimated  $\hat{Q}$

- Even though each estimate of the state-action values is unbiased, the estimate of  $\hat{\pi}$ 's value  $\hat{V}^{\hat{\pi}}$  can be biased:

$$\begin{aligned} \hat{V}^{\hat{\pi}} &= \mathbb{E} [\max (\hat{Q}(s, a_1), \hat{Q}(s, a_2))] \quad \text{expec. our random in rewards} \\ &\geq \max [\mathbb{E}(\hat{Q}(s, a_1)), \mathbb{E}(\hat{Q}(s, a_2))] \quad \text{Jensen's Ineq.} \\ &= \max [0, 0] = 0 = V^* \end{aligned}$$

<sup>39</sup>Example from Mannor, Simester, Sun and Tsitsiklis. Bias and Variance Approximation in Value Function Estimates. Management Science 2007

# Double Learning

- The greedy policy w.r.t. estimated  $Q$  values can yield a maximization bias during finite-sample learning
- Avoid using max of estimates as estimate of max of true values
- Instead split samples and use to create two independent unbiased estimates of  $Q_1(s_1, a_i)$  and  $Q_2(s_1, a_i) \forall a$ .

- Use one estimate to select max action:  $a^* = \arg \max_a Q_1(s_1, a)$

- Use other estimate to estimate value of  $a^*$ :  $Q_2(s, a^*)$

- Yields unbiased estimate:  $\mathbb{E}(Q_2(s, a^*)) = Q(s, a^*)$

- Why does this yield an unbiased estimate of the max state-action value?

$Q_2$  is an unbiased estimator of  $q_1$  &  $q_2$

- If acting online, can alternate samples used to update  $Q_1$  and  $Q_2$ , using the other to select the action chosen
- Next slides extend to full MDP case (with more than 1 state)

# Double Q-Learning

- 
- 1: Initialize  $Q_1(s, a)$  and  $Q_2(s, a), \forall s \in S, a \in A$   $t = 0$ , initial state  $s_t = s_0$
  - 2: **loop**
  - 3:   Select  $a_t$  using  $\epsilon$ -greedy  $\pi(s) = \arg \max_a Q_1(s_t, a) + Q_2(s_t, a)$
  - 4:   Observe  $(r_t, s_{t+1})$
  - 5:   **if** (with 0.5 probability) **then**
  - 6:      $Q_1(s_t, a_t) \leftarrow Q_1(s_t, a_t) + \alpha$
  - 7:   **else**
  - 8:      $Q_2(s_t, a_t) \leftarrow Q_2(s_t, a_t) + \alpha (r_t + \gamma \max_{a'} Q_1(s_{t+1}, a') - Q_2(s_t, a_t))$
  - 9:   **end if**
  - 10:    $t = t + 1$
  - 11: **end loop**
- 

*behavior policy*  
*Q learning off policy*

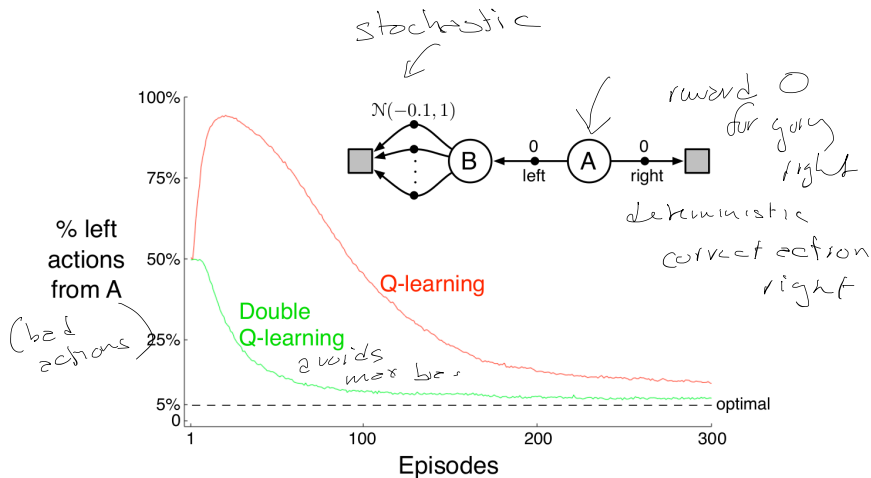
- Compared to Q-learning, how does this change the: memory requirements, computation requirements per step, amount of data required?

*doubles*

*same*

*much*

# Double Q-Learning (Figure 6.7 in Sutton and Barto 2018)



Due to the maximization bias, Q-learning spends much more time selecting suboptimal actions than double Q-learning.

# Table of Contents

- 1 Generalized Policy Iteration
- 2 Importance of Exploration
- 3 Monte Carlo Control
- 4 Temporal Difference Methods for Control
- 5 Maximization Bias

# Class Structure

- Last time: Policy evaluation with no knowledge of how the world works (MDP model not given)
- This time: Control (making decisions) without a model of how the world works
- **Next time: Value function approximation and Deep Q-learning**