

# CS292F StatRL Lecture 5

## RL Algorithms

Instructor: Yu-Xiang Wang

Spring 2021

UC Santa Barbara

# Homeworks and Project

- Homework 0 “due” tonight
  - Sticking to the schedule is highly recommended
- You should start forming project team / ideas now
  - Did you check out the list of papers that I sent?

# Recap: Lecture 4

- Miscellaneous on MDPs

- Advantage function and performance difference lemma
- Other types of MDPs: finite horizon, undiscounted, variable horizon...

- We started to talk about RL

- Model-based algorithms

- Model-free algorithms: Monte Carlo method, temporal difference methods.

$P$  with  $\uparrow$   
estimate  $P(s'/s, a)$   $\Rightarrow$   $\uparrow$  run  $VI$   
 $PI$

—  $\uparrow$   $\tau$   $t_0$  estimate  $V^{i_1}$   
—  $\uparrow$   $Q^*$   $t_0$  estimate  $Q^*$

# Recap: MDP planning with access to generative models

- Motivation:
  1. Solving MDP faster / approximately with randomized algs that sample
  2. Study sample complexity of RL with unknown transitions (without worrying about exploration)
- Algorithm of interest: Model-based plug-in estimator.
  - Sample all state-action pairs uniformly. Estimate the transition kernel.
  - Do VI / PI on the approximate MDP.

# Recap: "Monte Carlo" prediction / "Monte Carlo" control

- Sutton and Barto notations / terminologies:

- "Prediction"  $\Leftrightarrow$  "Policy evaluation" *fix  $\pi$  - compute  $V^\pi(\cdot)$ ,  $Q^\pi(\cdot)$*
- "Control"  $\Leftrightarrow$  "Policy optimization" *[how to come up with the next  $\pi$ ]*
- $Q$  and  $V$  are used to denote estimates / while  $q, v$  denotes the true value functions.
- 0 based indexing:  $S_0, A_0, R_1, S_1, A_1, R_2, \dots$

$q_\pi = Q^\pi$   
 $\uparrow$   
 in  $Q^\pi$  base  
 $\pi$  our policy

- Idea: Roll out trajectories and average them.

first visit MC:  $V_{(i)}^q(s) = \sum_{k=0}^{\infty} \gamma^k \cdot \text{index of first } S \cdot R_{k+1}^{(i)} = G_i$

$V(s) = \frac{1}{N} \sum_{i=1}^N G_i$        $\boxed{E[V^q(s)] = V^q(s)}$

first visit MC-control

$Q_{(i)}^\pi(s, a) = \sum_{k=0}^{\infty} \gamma^k \cdot \text{index of first } (S, a) \text{ pair} \cdot R_{k+1}^{(i)} = G_i$

$Q^\pi(s, a) = \frac{1}{N} \sum G_i$

# Recap: Temporal Difference Learning

- Monte Carlo  $V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$ ,

Issue:  $G_t$  can only be obtained after the entire episode!

- The idea of TD learning:

$$\mathbb{E}_\pi[G_t] = \mathbb{E}_\pi[R_t | S_t] + \gamma \mathbb{E}[V^\pi(S_{t+1})]$$

We only need one step before we can plug-in and estimate the RHS!

- TD-Policy evaluation

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

**Bootstrapping!**

# Recap: TD policy optimization (TD-control)

$\tau_1$  before  $S$  taken by  $A'$  taken by the previous  $\tau_1$   
 $S$  A R S A' R S A'

- SARSA (On-Policy TD-control)

Policy Eval

- Update the Q function by bootstrapping Bellman Equation

$E[Q(S,A)]$

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$$

- Choose the next A' using Q, e.g., eps-greedy.

Improvement

$$r(S,A) + \gamma E Q(S',A') - Q(S,A)$$

- Q-Learning (Off-policy TD-control)

- Update the Q function by bootstrapping Bellman Optimality Eq.

$$r(S,A) + \gamma E Q(S',A') - V^{\pi}(S)$$

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$

- Choose the next A' using Q, e.g., eps-greedy, or any other policy.

Remarks:

- These are **proven to converge** asymptotically.
- Much more data-efficient in practice, than MC.
- Regret analysis is still active area of research.

# Recap: Model-free vs Model-based RL algorithms

Model based

Model free

- Different function approximations

$\beta$

$V$   $Q$

- Different space efficiency

$O(S^2A)$

$O(SA)$

- Which one is more statistically efficient?
  - More or less equivalent in the tabular case.
  - Different challenges in their analysis.



# This lecture

- Variants / improvements of Sarsa and Q-learning
- TD-learning with function approximation
- Policy gradient methods

# Expected SARSA

Estimate  $G_t$  differently

SARSA'  $\bar{\pi}$

when  $\bar{\pi} = \text{argmax } Q$

$$\begin{aligned} Q(S_t, A_t) &\leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \mathbb{E}[Q(S_{t+1}, A_{t+1}) \mid S_{t+1}] - Q(S_t, A_t) \right] \\ &\leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \sum_a \underbrace{\pi(a|S_{t+1}) Q(S_{t+1}, a)}_{Q^*} - Q(S_t, A_t) \right], \end{aligned}$$

# Expected SARSA

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \mathbb{E}[Q(S_{t+1}, A_{t+1}) \mid S_{t+1}] - Q(S_t, A_t) \right]$$

$$\leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t) \right],$$

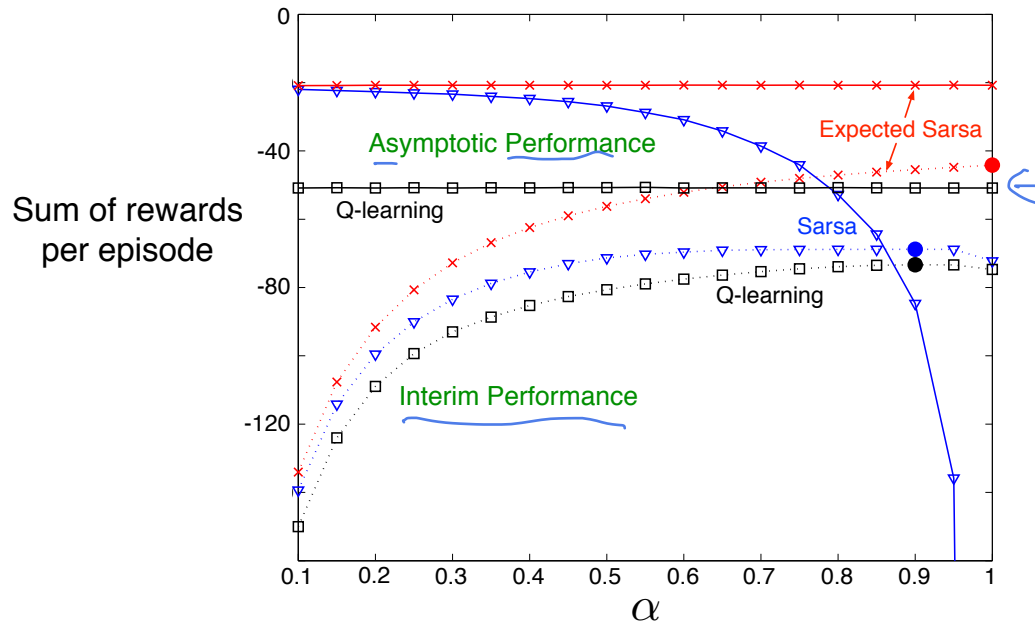


Figure 6.6: Interim and asymptotic performance of TD control methods on the cliff-walking task as a function of  $\alpha$ . All algorithms used an  $\epsilon$ -greedy policy with  $\epsilon = 0.1$ . Asymptotic performance is an average over 100,000 episodes whereas interim performance is an average over the first 100 episodes. These data are averages of over 50,000 and 10 runs for the interim and asymptotic cases respectively. The solid circles mark the best interim performance of each method. Adapted from van Seijen et al. (2009).

# Bias in Q-learning and double Q-learning

$$r(S_t, A_t) + \gamma \left( \max_a \mathbb{E}_{S_{t+1} \sim P(\cdot | S_t, A_t)} Q(S_{t+1}, a) \right) - Q(S_t, A_t)$$

estimate ↑

- Q-learning updates

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

$\gamma \left[ \mathbb{E} \max_a Q(S_{t+1}, a) \right]$ 
 $\mathbb{E} \max \geq \max \mathbb{E}$

- Double-Q-learning updates

- Keep track of two Q function estimates
- Toss a coin, if “head”:

$Q_1, Q_2$  are R.V.

$$Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \underbrace{Q_2(S_{t+1}, \arg \max_a Q_1(S_{t+1}, a))}_{a^* \perp Q_2} - Q_1(S_t, A_t) \right]$$

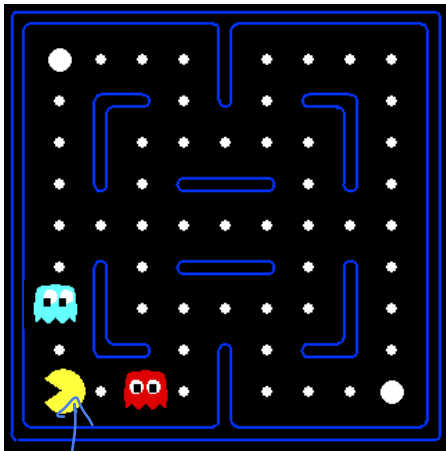
- If “tail”: update  $Q_2$

# The problem of large state-space is still there

- We need to represent and learn SA parameters in Q-learning and SARSA.
- S is often large
  - 9-puzzle, Tic-Tac-Toe:  $9! = 362,880$ ,  $S^2 = 1.3 \cdot 10^{11}$
  - PACMAN with 20 by 20 grid.  $S = O(2^{400})$ ,  $S^2 = O(2^{800})$
- $O(S)$  is not acceptable in some cases.
- Need to think of ways to “generalize”/share information across states.

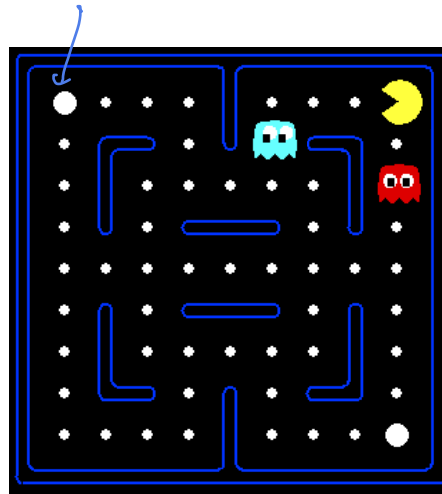
# Example: Pacman

Let's say we discover through experience that this state is bad:



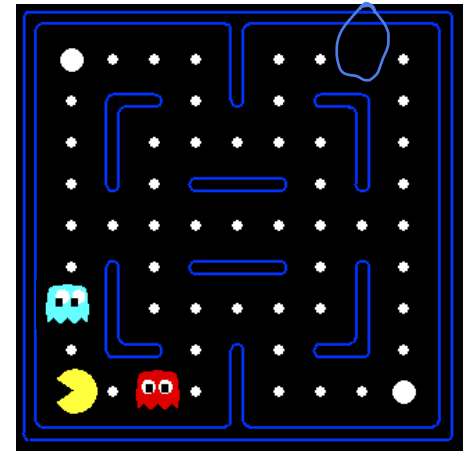
$S_1$

In naïve q-learning, we know nothing about this state:



$S_2$

Or even this one!



$S_3$

(From Dan Klein and Pieter Abbeel)

Video of Demo Q-Learning Pacman – Tiny – Watch All



# Video of Demo Q-Learning Pacman – Tiny – Silent Train



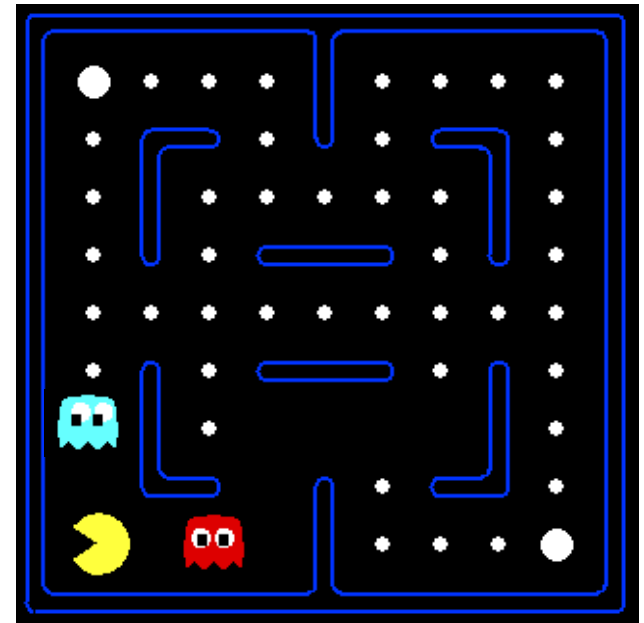


Video of Demo Q-Learning Pacman – Tricky –  
Watch All



# Why not use an evaluation function? A Feature-Based Representations

- Solution: describe a state using a vector of features (properties)
  - Features are functions from states to real numbers (often 0/1) that capture important properties of the state
  - Example features:
    - Distance to closest ghost
    - Distance to closest dot
    - Number of ghosts
    - $1 / (\text{dist to dot})^2$
    - Is Pacman in a tunnel? (0/1)
    - ..... etc.
    - Is it the exact state on this slide?
  - Can also describe a q-state (s, a) with features (e.g. action moves closer to food)



# Linear Value Functions

- Using a feature representation, we can write a q function (or value function) for any state using a few weights:
  - $V_{\mathbf{w}}(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$
  - $Q_{\mathbf{w}}(s,a) = w_1 f_1(s,a) + w_2 f_2(s,a) + \dots + w_n f_n(s,a)$
- Advantage: our experience is summed up in a few powerful numbers
- Disadvantage: states may share features but actually be very different in value!

# Updating a linear value function

- Original Q learning rule tries to reduce prediction error at  $s, a$ :

$$Q(s,a) \leftarrow Q(s,a) + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a) ]$$

- Instead, we update the weights to try to reduce the error at  $s, a$ :

$$\begin{aligned} w_i &\leftarrow w_i + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a) ] \partial Q_w(s,a) / \partial w_i \\ &= w_i + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a) ] f_i(s,a) \end{aligned}$$

# Updating a linear value function

- Original Q learning rule tries to reduce prediction error at  $s, a$ :

$$Q(s,a) \leftarrow Q(s,a) + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)]$$

- Instead, we update the weights to try to reduce the error at  $s, a$ :

$$\begin{aligned} w_i &\leftarrow w_i + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)] \partial Q_w(s,a) / \partial w_i \\ &= w_i + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)] f_i(s,a) \end{aligned}$$

- Qualitative justification:
  - Pleasant surprise: increase weights on positive features, decrease on negative ones
  - Unpleasant surprise: decrease weights on positive features, increase on negative ones

# PACMAN Q-Learning (Linear function approx.)



# Deriving the TD via incremental optimization that minimizes Bellman errors

- Mean Square Error and Mean Square Bellman error

$$\min_{w \in \mathbb{R}^d} \sum_s \mu(s) (V^{\pi}(s) - \hat{V}^{\pi}(s; w))^2$$

Monte Carlo: sample  $G_i^S$ ,  $V^{\pi}(s) = \frac{1}{N} \sum G_i^S$

$$\min_w \frac{1}{N} \sum_i \sum_s \mu(s) (G_i^S - \hat{V}^{\pi}(s; w))^2$$

$$\begin{aligned} w^T &= w - \alpha \nabla_w (G_i^S - \hat{V}^{\pi}(s; w))^2 \\ &= w + \alpha \underbrace{(-\hat{V}^{\pi}(s; w) + G_i^S)}_{\text{Monte Carlo TD estimate}} \cdot \nabla \hat{V}^{\pi}(s; w) \end{aligned}$$

*Bootstrapping!*

$$\min_{w \in \mathbb{R}^d} \sum_s \mu(s) (r(s) + \gamma E_{s' \sim p(s, \cdot)} \hat{V}^{\pi}(s'; w) - \hat{V}^{\pi}(s; w))^2$$

By sampling:  $S_1, A_1, R_2, S_2, A_2, R_3, \dots$

$$\min_{w \in \mathbb{R}^d} \frac{1}{T} \sum_{t=1}^T (R_{t+1} + \gamma \hat{V}^{\pi}(S_{t+1}; w) - \hat{V}^{\pi}(S_t; w))^2$$

$$\begin{aligned} w^T &= w - \alpha \nabla_w (R_{t+1} + \gamma \hat{V}^{\pi}(S_{t+1}; w) - \hat{V}^{\pi}(S_t; w))^2 \\ &\stackrel{\text{“Semi-gradient”}}{=} w + \alpha (R_{t+1} + \gamma \hat{V}^{\pi}(S_{t+1}; w) - \hat{V}^{\pi}(S_t; w)) \cdot \nabla \hat{V}^{\pi}(S_t; w) \end{aligned}$$

$f(S_t) \stackrel{\text{linear}}{=} \nabla \hat{V}^{\pi}(S_t; w)$

# How do we make sense of the semi-gradient method?

$$x_t := f(s_t) \in \mathbb{R}^d$$

$$x_{t+1} := f(s_{t+1}) \in \mathbb{R}^d$$

$$w_{t+1}^+ = w_t + \alpha [R_{t+1} + \gamma w_t^\top x_{t+1} - w_t^\top x_t] \cdot x_t$$

$$\frac{w^+ - w}{\alpha} = R_{t+1} \cdot x_t + x_t \cdot (\gamma x_{t+1} - x_t)^\top \cdot w_t$$

$$\dot{w} = R_{t+1} \cdot x_t + \underbrace{x_t \cdot (\gamma x_{t+1} - x_t)^\top}_{A} \cdot w_t$$

$$0 = \mathbb{E}[\dot{w} | w_t] = \mathbb{E}[\underbrace{R_{t+1}}_{\in \mathbb{R}^d} \cdot x_t] + \mathbb{E}[x_t \cdot (\gamma x_{t+1} - x_t)^\top] \cdot w_t$$

$$- A \in \mathbb{R}^{d \times d}$$

$$0 = b$$

$$Aw = b$$

$$-A w$$

$$\boxed{w^* = A^{-1} b} \leftarrow$$



# LSTD: Why doing incremental optimization when we can

$W^* = A^{-1}b$  try replace  $A$  with  $\hat{A}$ ,  $b$  with  $\hat{b}$

$$\hat{A}_t = \sum_{k=0}^{t-1} X_k (X_k - \gamma X_{k+1})^T + \varepsilon I$$

$$\hat{b}_t = \sum_{k=0}^{t-1} R_{k+1} \cdot X_k$$

$$\boxed{W_t = \hat{A}_t^{-1} \hat{b}_t} \quad O(d^3)$$

Rank 1 update of  $\hat{A}_{t-1}^{-1} \Rightarrow \hat{A}_t^{-1}$  Sherman-Morrison-Woodbury identity

$$\begin{aligned} \hat{A}_t^{-1} &= \left( \hat{A}_{t-1} + X_{t-1} (X_{t-1} - \gamma X_t)^T \right)^{-1} \\ &= \hat{A}_{t-1}^{-1} - \frac{\hat{A}_{t-1}^{-1} (X_{t-1} (X_{t-1} - \gamma X_t)^T) \hat{A}_{t-1}^{-1}}{1 + (X_{t-1} - \gamma X_t)^T \hat{A}_{t-1}^{-1} X_{t-1}} \end{aligned} \quad O(d^2)$$

# So far, in RL algorithms

- Model-based approaches
  - Estimate the MDP parameters.
  - Then use policy-iterations, value iterations.
- Monte Carlo methods:
  - estimating the rewards by empirical averages
- Temporal Difference methods:
  - Combine Monte Carlo methods with Dynamic Programming
- Linear function approximation in Q-learning
  - Similar to SGD
  - Learning heuristic function

\*Question: What is the policy class  $\Pi$  of interest in these methods?

# Remainder of the lecture

- Policy gradients methods
- Policy gradient theorem
- Extensions

# So far we talked about value function approximation, and an induced policy class.

S&B book  
AJKS book

- We can directly work with a parametric policy class.



- Examples:

1. Tabular setting: "Soft max" policy class,  
 $\theta \in \mathbb{R}^{S \times A}$

$$\pi_\theta(a|s) = \frac{\exp(\theta_{sa})}{\sum_{a'} \exp(\theta_{sa'})}$$

2. Linear feature setting, "log-linear policy class"

$$\pi_\theta(a|s) = \frac{\exp(\theta^T f(s,a))}{\sum_{a'} \exp(\theta^T f(s,a'))}$$

3. Nonlinear function class  
"neural policy class"

$$\pi_\theta(a|s) = \frac{\exp(f_\theta(s,a))}{\sum_{a'} \exp(f_\theta(s,a'))}$$

$$\sum_{a'} \exp(\theta^T f(s,a'))$$

$f_\theta$  is differentiable in  $\theta$

# Optimize policy using SGD

$$\min_{\theta} -V^{\pi_{\theta}}(u) \quad \square$$

$$\min_{\theta} \quad \square \quad \mathbb{E}_{\mathcal{S}_0 \sim \mathcal{D}_0} \left[ \sum_{t=0}^{\infty} \gamma^t R_{t+1} \right]$$

Stochastic optimization

$$\min_{\theta} -\frac{1}{N} \sum_{i=1}^N \sum R_{t+1}^{(i)}$$

other approximations with  $\theta$

$$\nabla_{\theta} V^{\pi_{\theta}}(u)$$

$$\theta^+ = \theta + \alpha \left[ \nabla_{\theta} V^{\pi_{\theta}}(u) \right]$$

$$\left[ \mathbb{E} \left[ \nabla_{\theta} V^{\pi_{\theta}}(u) \mid \theta \right] \right]$$

$$= \nabla_{\theta} V^{\pi_{\theta}}(u)$$

# How to estimate the gradient?

- Policy gradient theorem:

# Proof of Policy Gradient Theorem

# Proof of Policy Gradient Theorem



# REINFORCE Algorithm

REINFORCE, A Monte-Carlo Policy-Gradient Method (episodic)

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

Initialize policy parameter  $\boldsymbol{\theta} \in \mathbb{R}^{d'}$

Repeat forever:

    Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \boldsymbol{\theta})$

    For each step of the episode  $t = 0, \dots, T - 1$ :

$G \leftarrow$  return from step  $t$

$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \gamma^t G \nabla_{\boldsymbol{\theta}} \ln \pi(A_t|S_t, \boldsymbol{\theta})$

# Actor-Critic algorithm

- REINFORCE with a given baseline
- Actor-Critic: Learn the baseline and use the baseline for “bootstrapping”

# Next lecture

- Wrap up RL algorithms
- Exploration: Multi-armed bandits