

Lecture 4: Model Free Control and Function Approximation

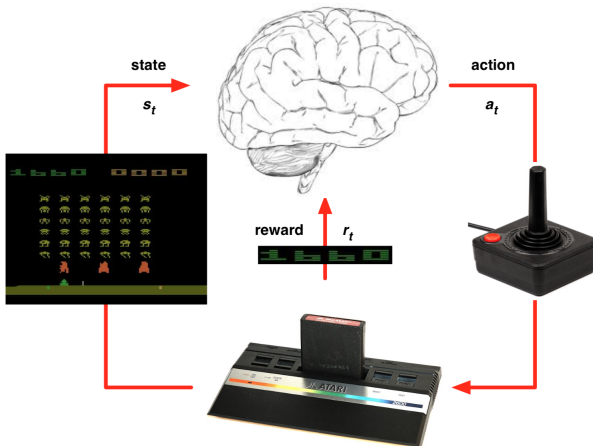
Emma Brunskill

CS234 Reinforcement Learning.

Winter 2026

- Structure and content drawn in part from David Silver's Lecture 5 and Lecture 6. For additional reading please see SB Sections 5.2-5.4, 6.4, 6.5, 6.7

Deep RL in Atari



Class Structure

- Last time: Policy evaluation with no knowledge of how the world works (MDP model not given)
- This time: first finish up policy evaluation when MDP model not given
- This time: Control (making decisions) without a model of how the world works
- Generalization – Value function approximation

Today's Lecture

- 1 Model Free Policy Evaluation in Tabular Settings
 - Batch MC and TD Policy Evaluation
- 2 Generalized Policy Improvement
 - Monte-Carlo Control with Tabular Representations
 - Greedy in the Limit of Infinite Exploration
 - Temporal Difference Methods for Control
- 3 Model Free Value Function Approximation
 - Policy Evaluation
 - Monte Carlo Policy Evaluation
 - Temporal Difference TD(0) Policy Evaluation
- 4 Control using Value Function Approximation
 - Control using General Value Function Approximators
 - Deep Q-Learning

Recap: MC, TD(0) and Certainty Equivalence Policy Evaluation

- Policy evaluation: Estimate $V^\pi(s)$ from executing π
- Trajectories τ : $(s, a \sim \pi(s), r, s', a' \sim \pi(s'), \dots)$ or tuples (s, a, r, s')
- MC: Given a full trajectory τ : $V^\pi(s) \leftarrow (1 - \alpha(s))V^\pi(s) + \alpha G_t(s)$
- TD(0): Given (s, a, r, s')
 $V^\pi(s) \leftarrow (1 - \alpha(s))V^\pi(s) + \alpha(s)(r + \gamma V^\pi(s'))$
- Certainty equivalence: Given a tuple (s, a, r, s') , update MLE dynamics model and reward model and then use policy evaluation methods to compute $V^\pi(s)$ for all s

Table of Contents

- 1 Model Free Policy Evaluation in Tabular Settings
 - Batch MC and TD Policy Evaluation
- 2 Generalized Policy Improvement
 - Monte-Carlo Control with Tabular Representations
 - Greedy in the Limit of Infinite Exploration
 - Temporal Difference Methods for Control
- 3 Model Free Value Function Approximation
 - Policy Evaluation
 - Monte Carlo Policy Evaluation
 - Temporal Difference TD(0) Policy Evaluation
- 4 Control using Value Function Approximation
 - Control using General Value Function Approximators
 - Deep Q-Learning

Table of Contents

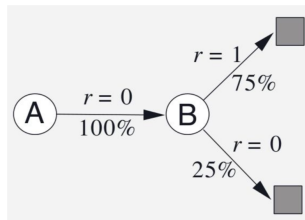
1 Model Free Policy Evaluation in Tabular Settings

- Batch MC and TD Policy Evaluation
- Monte-Carlo Control with Tabular Representations
- Greedy in the Limit of Infinite Exploration
- Temporal Difference Methods for Control
- Policy Evaluation
- Monte Carlo Policy Evaluation
- Temporal Difference TD(0) Policy Evaluation
- Control using General Value Function Approximators
- Deep Q-Learning

Batch MC and TD

- TD and MC methods shown use data once, then discard
- Batch (Offline) solution for finite dataset
 - Given set of K episodes
 - Repeatedly sample an episode from K
 - Apply MC or TD(0) to the sampled episode
- What do MC and TD(0) converge to?

AB Example: (Ex. 6.4, Sutton & Barto, 2018)



- Two states A, B with $\gamma = 1$
- Given 8 episodes of experience:
 - $A, 0, B, 0$
 - $B, 1$ (observed 6 times)
 - $B, 0$
- Imagine running TD updates over data infinite number of times
- $V(B) =$

AB Example: (Ex. 6.4, Sutton & Barto, 2018)

- TD Update: $V^\pi(s_t) = V^\pi(s_t) + \alpha(\underbrace{[r_t + \gamma V^\pi(s_{t+1})]}_{\text{TD target}} - V^\pi(s_t))$
- Two states A, B with $\gamma = 1$
- Given 8 episodes of experience:
 - $A, 0, B, 0$
 - $B, 1$ (observed 6 times)
 - $B, 0$
- Imagine run TD updates over data infinite number of times
- $V(B) = 0.75$ by TD or MC
- What about $V(A)$?

Check Your Understanding L3N3: AB Example: (Ex. 6.4, Sutton & Barto, 2018)

- TD Update: $V^\pi(s_t) = V^\pi(s_t) + \alpha(\underbrace{[r_t + \gamma V^\pi(s_{t+1})]}_{\text{TD target}} - V^\pi(s_t))$
- Two states A, B with $\gamma = 1$
- Given 8 episodes of experience:
 - $A, 0, B, 0$
 - $B, 1$ (observed 6 times)
 - $B, 0$
- Imagine run TD updates over data infinite number of times
- $V(B) = 0.75$ by TD or MC
- What about $V(A)$?
- Respond in Poll

Check Your Understanding L3N3: AB Example: (Ex. 6.4, Sutton & Barto, 2018)

- TD Update: $V^\pi(s_t) = V^\pi(s_t) + \alpha(\underbrace{[r_t + \gamma V^\pi(s_{t+1})]}_{\text{TD target}} - V^\pi(s_t))$
- Two states A, B with $\gamma = 1$
- Given 8 episodes of experience:
 - $A, 0, B, 0$
 - $B, 1$ (observed 6 times)
 - $B, 0$
- Imagine run TD updates over data infinite number of times
- $V(B) = 0.75$ by TD or MC
- What about $V(A)$?

Batch MC and TD: Convergence

Some Important Properties to Evaluate Model-free Policy Evaluation Algorithms

- Data efficiency & Computational efficiency
- In simple TD(0), use (s, a, r, s') once to update $V(s)$
 - $O(1)$ operation per update
 - In an episode of length L , $O(L)$
- In MC have to wait till episode finishes, then also $O(L)$
- MC can be more data efficient than simple TD
- But TD exploits Markov structure
 - If in Markov domain, leveraging this is helpful
- Dynamic programming with certainty equivalence also uses Markov structure

Summary: Policy Evaluation

Estimating the expected return of a particular policy if don't have access to true MDP models. Ex. evaluating average purchases per session of new product recommendation system

- Monte Carlo policy evaluation
 - Policy evaluation when we don't have a model of how the world works
 - Given on policy samples
 - Given off policy samples
- Temporal Difference (TD)
- Dynamic Programming with certainty equivalence
- *Understand what MC vs TD methods compute in batch evaluations
- Metrics / Qualities to evaluate and compare algorithms
 - Uses Markov assumption
 - Accuracy / MSE / bias / variance
 - Data efficiency
 - Computational efficiency

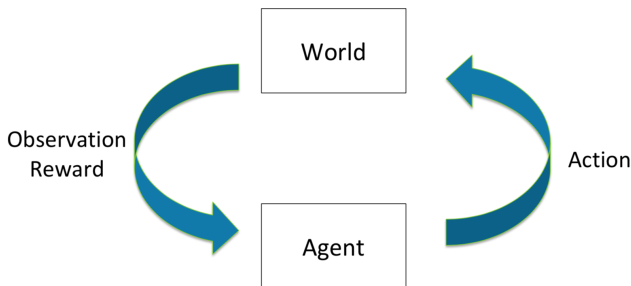
Table of Contents

- 1 Model Free Policy Evaluation in Tabular Settings
 - Batch MC and TD Policy Evaluation
- 2 Generalized Policy Improvement
 - Monte-Carlo Control with Tabular Representations
 - Greedy in the Limit of Infinite Exploration
 - Temporal Difference Methods for Control
- 3 Model Free Value Function Approximation
 - Policy Evaluation
 - Monte Carlo Policy Evaluation
 - Temporal Difference TD(0) Policy Evaluation
- 4 Control using Value Function Approximation
 - Control using General Value Function Approximators
 - Deep Q-Learning

Model-free Policy Iteration

- Initialize policy π
- Repeat:
 - Policy evaluation: compute Q^π
 - Policy improvement: update π given Q^π
- May need to modify policy evaluation:
 - If π 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

The Problem of Exploration



- Goal: Learn to select actions to maximize total expected future reward
- Problem: Can't learn about actions without trying them (need to *explore*)
- Problem: But if we try new actions, spending less time taking actions that our past experience suggests will yield high reward (need to *exploit* knowledge of domain to achieve high rewards)

ϵ -greedy Policies

- Simple idea to balance exploration and achieving rewards
- Let $|A|$ be the number of actions
- Then an ϵ -greedy policy w.r.t. a state-action value $Q(s, a)$ is $\pi(a|s) =$
 - $\arg \max_a Q(s, a)$, w. prob $1 - \epsilon + \frac{\epsilon}{|A|}$
 - $a' \neq \arg \max Q(s, a)$ w. prob $\frac{\epsilon}{|A|}$
- In words: select argmax action with probability $1 - \epsilon$, else select action uniformly at random

Policy Improvement with ϵ -greedy policies

- Recall we proved that policy iteration using given dynamics and reward models, was guaranteed to monotonically improve
- That proof assumed policy improvement output a deterministic policy
- Same property holds for ϵ -greedy policies

Monotonic ϵ -greedy Policy Improvement

Theorem

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

$$\begin{aligned} Q^{\pi_i}(s, \pi_{i+1}(s)) &= \sum_{a \in A} \pi_{i+1}(a|s) Q^{\pi_i}(s, a) \\ &= (\epsilon/|A|) \left[\sum_{a \in A} Q^{\pi_i}(s, a) \right] + (1 - \epsilon) \max_a Q^{\pi_i}(s, a) \end{aligned}$$

Table of Contents

- 1 Model Free Policy Evaluation in Tabular Settings
 - Batch MC and TD Policy Evaluation
- 2 Generalized Policy Improvement
 - Monte-Carlo Control with Tabular Representations
 - Greedy in the Limit of Infinite Exploration
 - Temporal Difference Methods for Control
- 3 Model Free Value Function Approximation
 - Policy Evaluation
 - Monte Carlo Policy Evaluation
 - Temporal Difference TD(0) Policy Evaluation
- 4 Control using Value Function Approximation
 - Control using General Value Function Approximators
 - Deep Q-Learning

Table of Contents

- Batch MC and TD Policy Evaluation
- Monte-Carlo Control with Tabular Representations
- Greedy in the Limit of Infinite Exploration
- Temporal Difference Methods for Control

3 Model Free Value Function Approximation

- Policy Evaluation
 - Monte Carlo Policy Evaluation
 - Temporal Difference TD(0) Policy Evaluation
- Control using General Value Function Approximators
- Deep Q-Learning

Recall Monte Carlo Policy Evaluation, Now for Q

```
1: Initialize  $Q(s, a) = 0, N(s, a) = 0 \forall (s, a), k = 1$ , Input  $\epsilon = 1, \pi$ 
2: loop
3:   Sample  $k$ -th episode  $(s_{k,1}, a_{k,1}, r_{k,1}, s_{k,2}, \dots, s_{k,T})$  given  $\pi$ 
3:   Compute  $G_{k,t} = r_{k,t} + \gamma r_{k,t+1} + \gamma^2 r_{k,t+2} + \dots \gamma^{T_i-1} r_{k,T_i} \forall t$ 
4:   for  $t = 1, \dots, T$  do
5:     if First visit to  $(s, a)$  in episode  $k$  then
6:        $N(s, a) = N(s, a) + 1$ 
7:        $Q(s_t, a_t) = Q(s_t, a_t) + \frac{1}{N(s, a)}(G_{k,t} - Q(s_t, a_t))$ 
8:     end if
9:   end for
10:   $k = k + 1$ 
11: end loop
```

Monte Carlo Online Control / On Policy Improvement

```
1: Initialize  $Q(s, a) = 0, N(s, a) = 0 \forall (s, a)$ , Set  $\epsilon = 1, k = 1$ 
2:  $\pi_k = \epsilon$ -greedy( $Q$ ) // Create initial  $\epsilon$ -greedy policy
3: loop
4:   Sample  $k$ -th episode  $(s_{k,1}, a_{k,1}, r_{k,1}, s_{k,2}, \dots, s_{k,T})$  given  $\pi_k$ 
4:    $G_{k,t} = r_{k,t} + \gamma r_{k,t+1} + \gamma^2 r_{k,t+2} + \dots + \gamma^{T-t} r_{k,T}$ 
5:   for  $t = 1, \dots, T$  do
6:     if First visit to  $(s, a)$  in episode  $k$  then
7:        $N(s, a) = N(s, a) + 1$ 
8:        $Q(s_t, a_t) = Q(s_t, a_t) + \frac{1}{N(s,a)}(G_{k,t} - Q(s_t, a_t))$ 
9:     end if
10:  end for
11:   $k = k + 1, \epsilon = 1/k$ 
12:   $\pi_k = \epsilon$ -greedy( $Q$ ) // Policy improvement
13: end loop
```

Optional Worked Example: MC for On Policy Control

- Mars rover with new actions:
 - $r(-, a_1) = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 +10]$, $r(-, a_2) = [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 +5]$, $\gamma = 1$.
- Assume current greedy $\pi(s) = a_1 \ \forall s$, $\epsilon=.5$. $Q(s, a) = 0$ for all (s, a)
- Sample trajectory from ϵ -greedy policy
- Trajectory = $(s_3, a_1, 0, s_2, a_2, 0, s_3, a_1, 0, s_2, a_2, 0, s_1, a_1, 1, \text{terminal})$
- First visit MC estimate of Q of each (s, a) pair?
- $Q^{\epsilon-\pi}(-, a_1) = [1 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$

After this trajectory (Select all)

- $Q^{\epsilon-\pi}(-, a_2) = [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0]$
- The new **greedy** policy would be: $\pi = [1 \ \text{tie} \ 1 \ \text{tie} \ \text{tie} \ \text{tie} \ \text{tie}]$
- The new **greedy** policy would be: $\pi = [1 \ 2 \ 1 \ \text{tie} \ \text{tie} \ \text{tie} \ \text{tie}]$
- If $\epsilon = 1/3$, prob of selecting a_1 in s_1 in the new ϵ -greedy policy is $1/9$.
- If $\epsilon = 1/3$, prob of selecting a_1 in s_1 in the new ϵ -greedy policy is $2/3$.
- If $\epsilon = 1/3$, prob of selecting a_1 in s_1 in the new ϵ -greedy policy is $5/6$.
- Not sure

Properties of MC control with ϵ -greedy policies

- Computational complexity?
- Converge to optimal Q^* function?
- Empirical performance?

Table of Contents

- Batch MC and TD Policy Evaluation
- Monte-Carlo Control with Tabular Representations
- Greedy in the Limit of Infinite Exploration
- Temporal Difference Methods for Control

3 Model Free Value Function Approximation

- Policy Evaluation
- Monte Carlo Policy Evaluation
- Temporal Difference TD(0) Policy Evaluation
- Control using General Value Function Approximators
- Deep Q-Learning

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 (policy used to act in the world) converges to greedy policy

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 (policy used to act in the world) converges to greedy policy
- A simple GLIE strategy is ϵ -greedy where ϵ is reduced to 0 with the following rate: $\epsilon_i = 1/i$

GLIE Monte-Carlo Control using Tabular Representations

Theorem

GLIE Monte-Carlo control converges to the optimal state-action value function $Q(s, a) \rightarrow Q^*(s, a)$

Table of Contents

- Batch MC and TD Policy Evaluation
- Monte-Carlo Control with Tabular Representations
- Greedy in the Limit of Infinite Exploration
- Temporal Difference Methods for Control

3 Model Free Value Function Approximation

- Policy Evaluation
- Monte Carlo Policy Evaluation
- Temporal Difference TD(0) Policy Evaluation
- Control using General Value Function Approximators
- Deep Q-Learning

Model-free Policy Iteration with TD Methods

- Initialize policy π
- Repeat:
 - Policy evaluation: compute Q^π using temporal difference updating with ϵ -greedy policy
 - Policy improvement: Same as Monte carlo policy improvement, set π to ϵ -greedy (Q^π)

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

Q-Learning: Learning the Optimal State-Action Value

- Q-learning
 - estimate the Q value of π^* while acting with another behavior policy π_b
- Key idea: Maintain Q estimates and bootstrap for best future value
- Q-learning:

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

Q-Learning with ϵ -greedy Exploration

-
- 1: Initialize $Q(s, a), \forall s \in S, a \in A$ $t = 0$, initial state $s_t = s_0$
 - 2: Set π_b to be ϵ -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: $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t))$
 - 7: $\pi(s_t) = \arg \max_a Q(s_t, a)$ w.prob $1 - \epsilon$, else random
 - 8: $t = t + 1$
 - 9: **end loop**
-

See optional worked example and optional understanding check at the end of the slides

Convergence Properties of Q-Learning

Theorem

Q-Learning 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 α_t satisfy the Robbins-Munro sequence such that

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

- For ex. $\alpha_t = \frac{1}{t}$ satisfies the above condition.

Properties of TD-Style Tabular Control with ϵ -greedy policies

- Result builds on stochastic approximation
- Relies on step sizes decreasing at the right rate
- Relies on Bellman backup contraction property
- Relies on bounded rewards and value function
- Note: other variants exist. SARSA (on-policy algorithm)
- SARSA

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

Table of Contents

- 1 Model Free Policy Evaluation in Tabular Settings
 - Batch MC and TD Policy Evaluation
- 2 Generalized Policy Improvement
 - Monte-Carlo Control with Tabular Representations
 - Greedy in the Limit of Infinite Exploration
 - Temporal Difference Methods for Control
- 3 Model Free Value Function Approximation
 - Policy Evaluation
 - Monte Carlo Policy Evaluation
 - Temporal Difference TD(0) Policy Evaluation
- 4 Control using Value Function Approximation
 - Control using General Value Function Approximators
 - Deep Q-Learning

Motivation for Function Approximation

- Avoid explicitly storing or learning the following for every single state and action
 - Dynamics or reward model
 - Value
 - State-action value
 - Policy
- Want more compact representation that generalizes across state or states and actions
 - Reduce memory needed to store $(P, R)/V/Q/\pi$
 - Reduce computation needed to compute $(P, R)/V/Q/\pi$
 - Reduce experience needed to find a good $(P, R)/V/Q/\pi$

State Action Value Function Approximation for Policy Evaluation with an Oracle

- First assume we could query any state s and action a and an oracle would return the true value for $Q^\pi(s, a)$
- Similar to supervised learning: assume given $((s, a), Q^\pi(s, a))$ pairs
- The objective is to find the best approximate representation of Q^π given a particular parameterized function $\hat{Q}(s, a; w)$

Stochastic Gradient Descent

- Goal: Find the parameter vector \mathbf{w} that minimizes the loss between a true value function $Q^\pi(s, a)$ and its approximation $\hat{Q}(s, a; \mathbf{w})$ as represented with a particular function class parameterized by \mathbf{w} .
- Generally use mean squared error and define the loss as

$$J(\mathbf{w}) = \mathbb{E}_\pi[(Q^\pi(s, a) - \hat{Q}(s, a; \mathbf{w}))^2]$$

- Can use gradient descent to find a local minimum

$$\Delta \mathbf{w} = -\frac{1}{2}\alpha \nabla_{\mathbf{w}} J(\mathbf{w})$$

- Stochastic gradient descent (SGD) uses a finite number of (often one) samples to compute an approximate gradient:
- Expected SGD is the same as the full gradient update

Stochastic Gradient Descent

- Goal: Find the parameter vector \mathbf{w} that minimizes the loss between a true value function $Q^\pi(s, a)$ and its approximation $\hat{Q}(s, a; \mathbf{w})$ as represented with a particular function class parameterized by \mathbf{w} .
- Generally use mean squared error and define the loss as

$$J(\mathbf{w}) = \mathbb{E}_\pi[(Q^\pi(s, a) - \hat{Q}(s, a; \mathbf{w}))^2]$$

- Can use gradient descent to find a local minimum

$$\Delta \mathbf{w} = -\frac{1}{2}\alpha \nabla_{\mathbf{w}} J(\mathbf{w})$$

- Stochastic gradient descent (SGD) uses a finite number of (often one) samples to compute an approximate gradient:

$$\begin{aligned}\nabla_{\mathbf{w}} J(\mathbf{w}) &= \nabla_{\mathbf{w}} E_\pi[Q^\pi(s, a) - \hat{Q}(s, a; \mathbf{w})]^2 \\ &= -2E_\pi[(Q^\pi(s, a) - \hat{Q}(s, a; \mathbf{w}))\nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w})]\end{aligned}$$

- Expected SGD is the same as the full gradient update

Table of Contents

- Batch MC and TD Policy Evaluation
 - Monte-Carlo Control with Tabular Representations
 - Greedy in the Limit of Infinite Exploration
 - Temporal Difference Methods for Control
 - Policy Evaluation
 - Monte Carlo Policy Evaluation
 - Temporal Difference TD(0) Policy Evaluation
-
- 4 Control using Value Function Approximation
 - Control using General Value Function Approximators
 - Deep Q-Learning

Model Free VFA Policy Evaluation

- No oracle to tell true $Q^\pi(s, a)$ for any state s and action a
- Use model-free state-action value function approximation

Model Free VFA Prediction / Policy Evaluation

- Recall model-free policy evaluation (Lecture 3)
 - Following a fixed policy π (or had access to prior data)
 - Goal is to estimate V^π and/or Q^π
- Maintained a lookup table to store estimates V^π and/or Q^π
- Updated these estimates after each episode (Monte Carlo methods) or after each step (TD methods)
- **Now: in value function approximation, change the estimate update step to include fitting the function approximator**

Table of Contents

- Batch MC and TD Policy Evaluation
- Monte-Carlo Control with Tabular Representations
- Greedy in the Limit of Infinite Exploration
- Temporal Difference Methods for Control
- Policy Evaluation
- Monte Carlo Policy Evaluation
- Temporal Difference TD(0) Policy Evaluation

4 Control using Value Function Approximation

- Control using General Value Function Approximators
- Deep Q-Learning

Monte Carlo Value Function Approximation

- Return G_t is an unbiased but noisy sample of the true expected return $Q^\pi(s_t, a_t)$
- Therefore can reduce MC VFA to doing supervised learning on a set of (state,action,return) pairs:
 $\langle (s_1, a_1), G_1 \rangle, \langle (s_2, a_2), G_2 \rangle, \dots, \langle (s_T, a_T), G_T \rangle$
 - Substitute G_t for the true $Q^\pi(s_t, a_t)$ when fit function approximator

MC Value Function Approximation for Policy Evaluation

```
1: Initialize  $\mathbf{w}$ ,  $k = 1$ 
2: loop
3:   Sample  $k$ -th episode  $(s_{k,1}, a_{k,1}, r_{k,1}, s_{k,2}, \dots, s_{k,L_k})$  given  $\pi$ 
4:   for  $t = 1, \dots, L_k$  do
5:     if First visit to  $(s, a)$  in episode  $k$  then
6:        $G_t(s, a) = \sum_{j=t}^{L_k} r_{k,j}$ 
7:        $\nabla_{\mathbf{w}} J(\mathbf{w}) = -2[G_t(s, a) - \hat{Q}(s_t, a_t; \mathbf{w})] \nabla_{\mathbf{w}} \hat{Q}(s_t, a_t; \mathbf{w})$  (Compute Gradient)
8:       Update weights  $\Delta \mathbf{w}$ 
9:     end if
10:  end for
11:   $k = k + 1$ 
12: end loop
```

Table of Contents

- Batch MC and TD Policy Evaluation
- Monte-Carlo Control with Tabular Representations
- Greedy in the Limit of Infinite Exploration
- Temporal Difference Methods for Control
- Policy Evaluation
- Monte Carlo Policy Evaluation
- Temporal Difference TD(0) Policy Evaluation

4 Control using Value Function Approximation

- Control using General Value Function Approximators
- Deep Q-Learning

Recall: Temporal Difference Learning w/ Lookup Table

- Uses bootstrapping and sampling to approximate V^π
- Updates $V^\pi(s)$ after each transition (s, a, r, s') :

$$V^\pi(s) = V^\pi(s) + \alpha(r + \gamma V^\pi(s') - V^\pi(s))$$

- Target is $r + \gamma V^\pi(s')$, a biased estimate of the true value $V^\pi(s)$
- Represent value for each state with a separate table entry

Temporal Difference TD(0) Learning with Value Function Approximation

- Uses bootstrapping and sampling to approximate true V^π
- Updates estimate $V^\pi(s)$ after each transition (s, a, r, s') :

$$V^\pi(s) = V^\pi(s) + \alpha(r + \gamma V^\pi(s') - V^\pi(s))$$

- Target is $r + \gamma V^\pi(s')$, a biased estimate of the true value $V^\pi(s)$
- In value function approximation, target is $r + \gamma \hat{V}^\pi(s'; \mathbf{w})$, a biased and approximated estimate of the true value $V^\pi(s)$
- 3 forms of approximation:
 - 1 Sampling
 - 2 Bootstrapping
 - 3 Value function approximation

Temporal Difference TD(0) Learning with Value Function Approximation

- In value function approximation, target is $r + \gamma \hat{V}^\pi(s'; \mathbf{w})$, a biased and approximated estimate of the true value $V^\pi(s)$
- Can reduce doing TD(0) learning with value function approximation to supervised learning on a set of data pairs:
 - $\langle s_1, r_1 + \gamma \hat{V}^\pi(s_2; \mathbf{w}) \rangle, \langle s_2, r_2 + \gamma \hat{V}^\pi(s_3; \mathbf{w}) \rangle, \dots$
- Find weights to minimize mean squared error

$$J(\mathbf{w}) = \mathbb{E}_\pi[(r_j + \gamma \hat{V}^\pi(s_{j+1}, \mathbf{w}) - \hat{V}(s_j; \mathbf{w}))^2]$$

- Use stochastic gradient descent, as in MC methods

TD(0) Value Function Approximation for Policy Evaluation

-
- 1: Initialize \mathbf{w}, \mathbf{s}
 - 2: **loop**
 - 3: Given s sample $a \sim \pi(s), r(s, a), s' \sim p(s'|s, a)$
 - 4: $\nabla_{\mathbf{w}} J(\mathbf{w}) = -2[r + \gamma \hat{V}(s'; \mathbf{w}) - \hat{V}(s; \mathbf{w})] \nabla_{\mathbf{w}} \hat{V}(s; \mathbf{w})$
 - 5: Update weights $\Delta \mathbf{w}$
 - 6: **if** s' is not a terminal state **then**
 - 7: Set $s = s'$
 - 8: **else**
 - 9: Restart episode, sample initial state s
 - 10: **end if**
 - 11: **end loop**
-

Table of Contents

- 1 Model Free Policy Evaluation in Tabular Settings
 - Batch MC and TD Policy Evaluation
- 2 Generalized Policy Improvement
 - Monte-Carlo Control with Tabular Representations
 - Greedy in the Limit of Infinite Exploration
 - Temporal Difference Methods for Control
- 3 Model Free Value Function Approximation
 - Policy Evaluation
 - Monte Carlo Policy Evaluation
 - Temporal Difference TD(0) Policy Evaluation
- 4 Control using Value Function Approximation
 - Control using General Value Function Approximators
 - Deep Q-Learning

Table of Contents

- Batch MC and TD Policy Evaluation
- Monte-Carlo Control with Tabular Representations
- Greedy in the Limit of Infinite Exploration
- Temporal Difference Methods for Control
- Policy Evaluation
- Monte Carlo Policy Evaluation
- Temporal Difference TD(0) Policy Evaluation
- Control using General Value Function Approximators
- Deep Q-Learning

Control using Value Function Approximation

- Use value function approximation to represent state-action values
 $\hat{Q}^{\pi}(s, a; \mathbf{w}) \approx Q^{\pi}$
- Interleave
 - Approximate policy evaluation using value function approximation
 - Perform ϵ -greedy policy improvement
- Can be unstable. Generally involves intersection of the following:
 - Function approximation
 - Bootstrapping
 - **Off-policy learning**

Action-Value Function Approximation with an Oracle

- $\hat{Q}^\pi(s, a; \mathbf{w}) \approx Q^\pi$
- Minimize the mean-squared error between the true action-value function $Q^\pi(s, a)$ and the approximate action-value function:

$$J(\mathbf{w}) = \mathbb{E}_\pi[(Q^\pi(s, a) - \hat{Q}^\pi(s, a; \mathbf{w}))^2]$$

- Use stochastic gradient descent to find a local minimum

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = -2\mathbb{E}[(Q^\pi(s, a) - \hat{Q}^\pi(s, a; \mathbf{w}))\nabla_{\mathbf{w}} \hat{Q}^\pi(s, a; \mathbf{w})]$$

- Stochastic gradient descent (SGD) samples the gradient

Incremental Model-Free Control Approaches

- Similar to policy evaluation, true state-action value function for a state is unknown and so substitute a target value for true $Q(s_t, a_t)$

$$\Delta \mathbf{w} = \alpha(Q(s_t, a_t) - \hat{Q}(s_t, a_t; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s_t, a_t; \mathbf{w})$$

- In Monte Carlo methods, use a return G_t as a substitute target

$$\Delta \mathbf{w} = \alpha(G_t - \hat{Q}(s_t, a_t; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s_t, a_t; \mathbf{w})$$

- SARSA: Use TD target $r + \gamma \hat{Q}(s', a'; \mathbf{w})$ which leverages the current function approximation value

$$\Delta \mathbf{w} = \alpha(r + \gamma \hat{Q}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w})$$

- Q-learning: Uses related TD target $r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w})$

$$\Delta \mathbf{w} = \alpha(r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w})$$

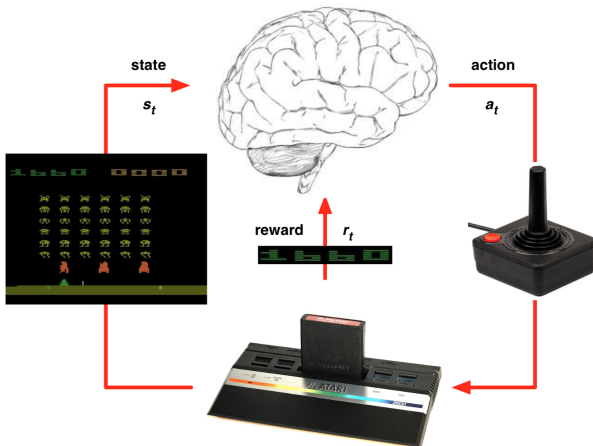
"Deadly Triad" which Can Cause Instability

- Informally, updates involve doing an (approximate) Bellman backup followed by best trying to fit underlying value function to a particular feature representation
- Bellman operators are contractions, but value function approximation fitting can be an expansion
 - To learn more, see Baird example in Sutton and Barto 2018
- "Deadly Triad" can lead to oscillations or lack of convergence
 - Bootstrapping
 - Function Approximation
 - Off policy learning (e.g. Q-learning)

Table of Contents

- Batch MC and TD Policy Evaluation
- Monte-Carlo Control with Tabular Representations
- Greedy in the Limit of Infinite Exploration
- Temporal Difference Methods for Control
- Policy Evaluation
- Monte Carlo Policy Evaluation
- Temporal Difference TD(0) Policy Evaluation
- Control using General Value Function Approximators
- Deep Q-Learning

Using these ideas to do Deep RL in Atari



Q-Learning with Neural Networks

- Q-learning converges to optimal $Q^*(s, a)$ using tabular representation
- In value function approximation Q-learning minimizes MSE loss by stochastic gradient descent using a target Q estimate instead of true Q
- But Q-learning with VFA can diverge
- Two of the issues causing problems:
 - Correlations between samples
 - Non-stationary targets
- Deep Q-learning (DQN) addresses these challenges by using
 - Experience replay
 - Fixed Q-targets

DQNs: Experience Replay

- To help remove correlations, store dataset (called a **replay buffer**) \mathcal{D} from prior experience

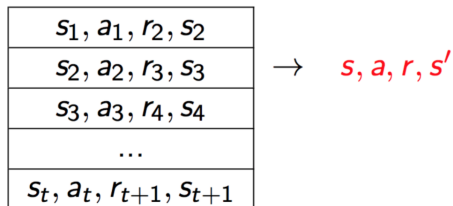
| | |
|------------------------------|---------------------------|
| s_1, a_1, r_2, s_2 | $\rightarrow s, a, r, s'$ |
| s_2, a_2, r_3, s_3 | |
| s_3, a_3, r_4, s_4 | |
| ... | |
| $s_t, a_t, r_{t+1}, s_{t+1}$ | |

- To perform experience replay, repeat the following:
 - $(s, a, r, s') \sim \mathcal{D}$: sample an experience tuple from the dataset
 - Compute the target value for the sampled s : $r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w})$
 - Use stochastic gradient descent to update the network weights

$$\Delta \mathbf{w} = \alpha (r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w})$$

DQNs: Experience Replay

- To help remove correlations, store dataset \mathcal{D} from prior experience



- To perform experience replay, repeat the following:
 - $(s, a, r, s') \sim \mathcal{D}$: sample an experience tuple from the dataset
 - Compute the target value for the sampled s : $r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w})$
 - Use stochastic gradient descent to update the network weights

$$\Delta \mathbf{w} = \alpha (r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w})$$

- **Uses target as a scalar, but function weights will get updated on the next round, changing the target value**

DQNs: Fixed Q-Targets

- To help improve stability, fix the **target weights** used in the target calculation for multiple updates
- Target network uses a different set of weights than the weights being updated
- Let parameters \mathbf{w}^- be the set of weights used in the target, and \mathbf{w} be the weights that are being updated
- Slight change to computation of target value:
 - $(s, a, r, s') \sim \mathcal{D}$: sample an experience tuple from the dataset
 - Compute the target value for the sampled s : $r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w}^-)$
 - Use stochastic gradient descent to update the network weights

$$\Delta \mathbf{w} = \alpha (r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w}^-) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w})$$

DQN Pseudocode

```
1: Input  $C, \alpha, D = \{\}$ , Initialize  $\mathbf{w}, \mathbf{w}^- = \mathbf{w}, t = 0$ 
2: Get initial state  $s_0$ 
3: loop
4:   Sample action  $a_t$  given  $\epsilon$ -greedy policy for current  $\hat{Q}(s_t, a; \mathbf{w})$ 
5:   Observe reward  $r_t$  and next state  $s_{t+1}$ 
6:   Store transition  $(s_t, a_t, r_t, s_{t+1})$  in replay buffer  $D$ 
7:   Sample random minibatch of tuples  $(s_i, a_i, r_i, s_{i+1})$  from  $D$ 
8:   for  $j$  in minibatch do
9:     if episode terminated at step  $i + 1$  then
10:       $y_i = r_i$ 
11:     else
12:       $y_i = r_i + \gamma \max_{a'} \hat{Q}(s_{i+1}, a'; \mathbf{w}^-)$ 
13:     end if
14:     Do gradient descent step on  $(y_i - \hat{Q}(s_i, a_i; \mathbf{w}))^2$  for parameters  $\mathbf{w}$ :  $\Delta \mathbf{w} = \alpha(y_i - \hat{Q}(s_i, a_i; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s_i, a_i; \mathbf{w})$ 
15:   end for
16:    $t = t + 1$ 
17:   if  $\text{mod}(t, C) == 0$  then
18:      $\mathbf{w}^- \leftarrow \mathbf{w}$ 
19:   end if
20: end loop
```

Note there are several hyperparameters and algorithm choices. One needs to choose the neural network architecture, the learning rate, and how often to update the target network. Often a fixed size replay buffer is used for experience replay, which introduces a parameter to control the size, and the need to decide how to populate it.

Check Your Understanding L4N3: Fixed Targets

- In DQN we compute the target value for the sampled (s, a, r, s') using a separate set of target weights: $r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w}^-)$
- Select all that are true
- This doubles the computation time compared to a method that does not have a separate set of weights
- This doubles the memory requirements compared to a method that does not have a separate set of weights
- Not sure

Check Your Understanding L4N3: Fixed Targets.

Solutions

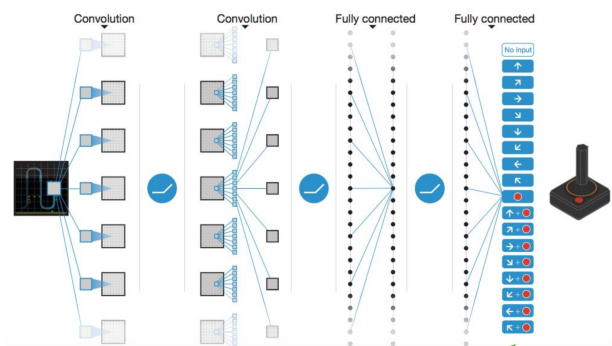
- In DQN we compute the target value for the sampled (s, a, r, s') using a separate set of target weights: $r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w}^-)$
- Select all that are true
- This doubles the computation time compared to a method that does not have a separate set of weights
- This doubles the memory requirements compared to a method that does not have a separate set of weights
- Not sure

DQNs Summary

- DQN uses experience replay and fixed Q-targets
- Store transition $(s_t, a_t, r_{t+1}, s_{t+1})$ in replay memory \mathcal{D}
- Sample random mini-batch of transitions (s, a, r, s') from \mathcal{D}
- Compute Q-learning targets w.r.t. old, fixed parameters \mathbf{w}^-
- Optimizes MSE between Q-network and Q-learning targets
- Uses stochastic gradient descent

DQNs in Atari

- End-to-end learning of values $Q(s, a)$ from pixels s
- Input state s is stack of raw pixels from last 4 frames
- Output is $Q(s, a)$ for 18 joystick/button positions
- Reward is change in score for that step
- Used a deep neural network with CNN
- Network architecture and hyperparameters fixed across all games



1 network, outputs Q value for each action

Figure: Human-level control through deep reinforcement learning, Mnih et al, 2015

DQN Results in Atari

Figure: Human-level control through deep reinforcement learning, Mnih et al, 2015

Which Aspects of DQN were Important for Success?

| Game | Linear | Deep Network |
|----------------|--------|--------------|
| Breakout | 3 | 3 |
| Enduro | 62 | 29 |
| River Raid | 2345 | 1453 |
| Seaquest | 656 | 275 |
| Space Invaders | 301 | 302 |

Note: just using a deep NN actually hurt performance sometimes!

Which Aspects of DQN were Important for Success?

| Game | Linear | Deep Network | DQN w/ fixed Q |
|----------------|--------|--------------|-------------------|
| Breakout | 3 | 3 | 10 |
| Enduro | 62 | 29 | 141 |
| River Raid | 2345 | 1453 | 2868 |
| Seaquest | 656 | 275 | 1003 |
| Space Invaders | 301 | 302 | 373 |

Which Aspects of DQN were Important for Success?

| Game | Linear | Deep Network | DQN w/ fixed Q | DQN w/ replay | DQN w/replay and fixed Q |
|----------------|--------|--------------|----------------|---------------|--------------------------|
| Breakout | 3 | 3 | 10 | 241 | 317 |
| Enduro | 62 | 29 | 141 | 831 | 1006 |
| River Raid | 2345 | 1453 | 2868 | 4102 | 7447 |
| Seaquest | 656 | 275 | 1003 | 823 | 2894 |
| Space Invaders | 301 | 302 | 373 | 826 | 1089 |

- Replay is **hugely** important
- Why? Beyond helping with correlation between samples, what does replaying do?

- Success in Atari has led to huge excitement in using deep neural networks to do value function approximation in RL
- Some immediate improvements (many others!)
 - **Double DQN** (Deep Reinforcement Learning with Double Q-Learning, Van Hasselt et al, AAAI 2016)
 - Prioritized Replay (Prioritized Experience Replay, Schaul et al, ICLR 2016)
 - Dueling DQN (best paper ICML 2016) (Dueling Network Architectures for Deep Reinforcement Learning, Wang et al, ICML 2016)

What You Should Understand

- Be able to implement TD(0) and MC on policy evaluation
- Be able to implement Q-learning and SARSA and MC control algorithms
- List the 3 issues that can cause instability and describe the problems qualitatively: function approximation, bootstrapping and off-policy learning
- Know some of the key features in DQN that were critical (experience replay, fixed targets)

Class Structure

- Last time and start of this time: Model-free reinforcement learning with function approximation
- Next time: Policy gradients

Monotonic ϵ -greedy Policy Improvement

Theorem

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

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

SARSA Initialization Conceptual Question

- Mars rover with new actions:
 - $r(-, a_1) = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 + 10]$, $r(-, a_2) = [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 + 5]$, $\gamma = 1$.
- Initialize $\epsilon = 1/k$, $k = 1$, and $\alpha = 0.5$, $Q(-, a_1) = r(-, a_1)$,
 $Q(-, a_2) = r(-, a_2)$
- SARSA: $(s_6, a_1, 0, s_7, a_2, 5, s_7)$.
- Does how Q is initialized matter (initially? asymptotically)?

Optional Worked Example: MC for On Policy Control Solution

- Mars rover with new actions:
 - $r(-, a_1) = [1 \ 0 \ 0 \ 0 \ 0 \ 0 +10]$, $r(-, a_2) = [0 \ 0 \ 0 \ 0 \ 0 \ 0 +5]$, $\gamma = 1$.
- Assume current greedy $\pi(s) = a_1 \ \forall s$, $\epsilon = .5$. $Q(s, a) = 0$ for all (s, a)
- Sample trajectory from ϵ -greedy policy
- Trajectory = $(s_3, a_1, 0, s_2, a_2, 0, s_3, a_1, 0, s_2, a_2, 0, s_1, a_1, 1, \text{terminal})$
- First visit MC estimate of Q of each (s, a) pair?
- $Q^{\epsilon-\pi}(-, a_1) = [1 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0]$

After this trajectory:

- $Q^{\epsilon-\pi}(-, a_2) = [0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0]$
- The new **greedy** policy would be: $\pi = [1 \ 2 \ 1 \ \text{tie} \ \text{tie} \ \text{tie} \ \text{tie}]$
- If $\epsilon = 1/3$, prob of selecting a_1 in s_1 in the new ϵ -greedy policy is $5/6$.

Optional Worked Example SARSA for Mars Rover

```
1: Set initial  $\epsilon$ -greedy policy  $\pi$ ,  $t = 0$ , initial state  $s_t = s_0$ 
2: Take  $a_t \sim \pi(s_t)$  // Sample action from policy
3: Observe  $(r_t, s_{t+1})$ 
4: loop
5:   Take action  $a_{t+1} \sim \pi(s_{t+1})$ 
6:   Observe  $(r_{t+1}, s_{t+2})$ 
7:    $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$ 
8:    $\pi(s_t) = \arg \max_a Q(s_t, a)$  w.prob  $1 - \epsilon$ , else random
9:    $t = t + 1$ 
10: end loop
```

- Initialize $\epsilon = 1/k$, $k = 1$, and $\alpha = 0.5$, $Q(-, a_1) = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ +10]$, $Q(-, a_2) = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ +5]$, $\gamma = 1$
- Assume starting state is s_6 and sample a_1

Worked Example: SARSA for Mars Rover

-
- 1: Set initial ϵ -greedy policy π , $t = 0$, initial state $s_t = s_0$
 - 2: Take $a_t \sim \pi(s_t)$ // Sample action from policy
 - 3: Observe (r_t, s_{t+1})
 - 4: **loop**
 - 5: Take action $a_{t+1} \sim \pi(s_{t+1})$
 - 6: Observe (r_{t+1}, s_{t+2})
 - 7: $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$
 - 8: $\pi(s_t) = \arg \max_a Q(s_t, a)$ w.prob $1 - \epsilon$, else random
 - 9: $t = t + 1$
 - 10: **end loop**
-
- Initialize $\epsilon = 1/k$, $k = 1$, and $\alpha = 0.5$, $Q(-, a_1) = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ +10]$, $Q(-, a_2) = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ +5]$, $\gamma = 1$
 - Assume starting state is s_6 and sample a_1

Worked Example: SARSA for Mars Rover

-
- 1: Set initial ϵ -greedy policy π , $t = 0$, initial state $s_t = s_0$
 - 2: Take $a_t \sim \pi(s_t)$ // Sample action from policy
 - 3: Observe (r_t, s_{t+1})
 - 4: **loop**
 - 5: Take action $a_{t+1} \sim \pi(s_{t+1})$
 - 6: Observe (r_{t+1}, s_{t+2})
 - 7: $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$
 - 8: $\pi(s_t) = \arg \max_a Q(s_t, a)$ w.prob $1 - \epsilon$, else random
 - 9: $t = t + 1$
 - 10: **end loop**
-

- Initialize $\epsilon = 1/k$, $k = 1$, and $\alpha = 0.5$, $Q(-, a_1) = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ +10]$, $Q(-, a_2) = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ +5]$, $\gamma = 1$
- Tuple: $(s_6, a_1, 0, s_7, a_2, 5, s_7)$.
- $Q(s_6, a_1) = .5 * 0 + .5 * (0 + \gamma Q(s_7, a_2)) = 2.5$

Worked Example: ϵ -greedy Q-Learning Mars

-
- 1: Initialize $Q(s, a), \forall s \in S, a \in A$ $t = 0$, initial state $s_t = s_0$
 - 2: Set π_b to be ϵ -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: $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t))$
 - 7: $\pi(s_t) = \arg \max_a Q(s_t, a)$ w.prob $1 - \epsilon$, else random
 - 8: $t = t + 1$
 - 9: **end loop**
-

- Initialize $\epsilon = 1/k$, $k = 1$, and $\alpha = 0.5$, $Q(-, a_1) = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ +10]$, $Q(-, a_2) = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ +5]$, $\gamma = 1$
- Like in SARSA example, start in s_6 and take a_1 .

Worked Example: ϵ -greedy Q-Learning Mars

-
- 1: Initialize $Q(s, a), \forall s \in S, a \in A$ $t = 0$, initial state $s_t = s_0$
 - 2: Set π_b to be ϵ -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: $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t))$
 - 7: $\pi(s_t) = \arg \max_a Q(s_t, a)$ w.prob $1 - \epsilon$, else random
 - 8: $t = t + 1$
 - 9: **end loop**
-

- Initialize $\epsilon = 1/k$, $k = 1$, and $\alpha = 0.5$, $Q(-, a_1) = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ +10]$, $Q(-, a_2) = [1 \ 0 \ 0 \ 0 \ 0 \ 0 \ +5]$, $\gamma = 1$
- Tuple: $(s_6, a_1, 0, s_7)$.
- $Q(s_6, a_1) = 0 + .5 * (0 + \gamma \max_{a'} Q(s_7, a') - 0) = .5 * 10 = 5$
- Recall that in the SARSA update we saw $Q(s_6, a_1) = 2.5$ because we used the actual action taken at s_7 instead of the max
- Does how Q is initialized matter (initially? asymptotically?)?

Optional Check Your Understanding L4: SARSA and Q-Learning

- SARSA: $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$
- Q-Learning:
 $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t))$

Select all that are true

- 1 Both SARSA and Q-learning may update their policy after every step
- 2 If $\epsilon = 0$ for all time steps, and Q is initialized randomly, a SARSA Q state update will be the same as a Q-learning Q state update
- 3 Not sure

Optional Check Your Understanding SARSA and Q-Learning Solutions

- SARSA: $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$
- Q-Learning:
 $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha(r_t + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t))$

Select all that are true

- 1 Both SARSA and Q-learning may update their policy after every step
- 2 If $\epsilon = 0$ for all time steps, and Q is initialized randomly, a SARSA Q state update will be the same as a Q-learning Q state update
- 3 Not sure