Generalized Policy Iteration & Evaluating Value Fcts without dynamics: Monte-Carlo, Temporal Difference Learning



Agent and environment interact at discrete time steps: t = 0, 1, 2, 3, ...Agent observes state at step t: $S_t \in S$ produces action at step t: $A_t \in \mathcal{A}(S_t)$ gets resulting reward: $R_{t+1} \in \mathcal{R} \subset \mathbb{R}$ and resulting next state: $S_{t+1} \in S^+$

Recall: Return

Agent wants to maximize it's return:

$$G_{t} = R_{t+1} + \gamma R_{t+2} + \gamma^{2} R_{t+3} + L = \sum_{k=0}^{\infty} \gamma^{k} R_{t+k+1},$$

. . .

where $\gamma, 0 \le \gamma \le 1$, is the **discount rate**.

shortsighted $0 \leftarrow \gamma \rightarrow 1$ farsighted

4 value functions

	state values	action values
prediction	v_{π}	q_{π}
control	v_*	q_*

- All theoretical objects, expected values
- Distinct from their estimates: $V_t(s) = Q_t(s,a)$

Today: Algorithms to Estimate v, q

Generalized Policy Iteration (improving a policy)

□ MC: Monte-Carlo

TD: Temporal Difference Learning

Bellman Equations

Beliman Eqn:

$$v_{\pi}(s) = \sum_{a} \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma v_{\pi}(s')]$$
Beliman Optimality Eqn:

$$v_{*}(s) = \max_{a} \sum_{s',r} p(s',r|s,a) [r + \gamma v_{*}(s')]$$

*Also as many equations as unknowns (non-linear, this time though).

Policy Iteration – One array version (+ policy)

- 1. Initialization $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in S$
- 2. Policy Evaluation

Repeat

 $\Delta \leftarrow 0, v \leftarrow V$ For each $s \in S$:

$$V(s) \leftarrow \sum_{s',r} p(s',r \,|\, s,\pi(s)) \big[r + \gamma v(s') \big]$$

until $\Delta < \theta$ (a small positive number)

3. Policy Improvement policy-stable \leftarrow true For each $s \in S$: $a \leftarrow \pi(s)$ $\pi(s) \leftarrow \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$ If $a \neq \pi(s)$, then policy-stable \leftarrow false If policy-stable, then stop and return V and π ; else go to 2

Generalized Policy Iteration

Generalized Policy Iteration (GPI): any interaction of policy evaluation and policy improvement, independent of their granularity.



A geometric metaphor for convergence of GPI:



Value Iteration

Recall the **full policy-evaluation backup**:

$$v_{k+1}(s) = \sum_{a} \pi(a|s) \sum_{s',r} p(s',r|s,a) \left[r + \gamma v_k(s')\right] \qquad \forall s \in \mathcal{S}$$

Here is the **full value-iteration backup**:

$$v_{k+1}(s) = \max_{a} \sum_{s',r} p(s',r|s,a) \left[r + \gamma v_k(s')\right] \qquad \forall s \in \mathcal{S}$$

Initialize array V arbitrarily (e.g., V(s) = 0 for all $s \in S^+$)

Repeat $\Delta \leftarrow 0, v \leftarrow V$ For each $s \in S$: $V(s) \leftarrow \max_{a} \sum_{s',r} p(s',r|s,a) [r + \gamma v(s')]$ $\Delta \leftarrow \max(\Delta, |v(s) - V(s)|)$ until $\Delta < \theta$ (a small positive number)

Output a deterministic policy, π , such that $\pi(s) = \arg \max_a \sum_{s',r} p(s', r|s, a) [r + \gamma V(s')]$

Gambler's Problem

- □ Gambler can repeatedly bet \$ on a coin flip
- Heads he wins his stake, tails he loses it
- □ Initial capital $\in \{\$1, \$2, \dots \$99\}$
- Gambler wins if his capital becomes \$100 loses if it becomes \$0
- **Coin is unfair**
 - Heads (gambler wins) with probability p = .4
- **I** States, Actions, Rewards? Discounting?

Gambler's Problem Solution



Gambler's Problem Solution



Herd Management

- □ You are a consultant to a farmer managing a herd of cows
- Herd consists of 5 kinds of cows:
 - Young
 - Milking
 - Breeding
 - Old
 - Sick
- Number of each kind is the State
- Number sold of each kind is the Action
- **Cows transition from one kind to another**
- ☐ Young cows can be born

Asynchronous DP

- All the DP methods described so far require exhaustive sweeps of the entire state set.
- Asynchronous DP does not use sweeps. Instead it works like this:
 - Repeat until convergence criterion is met:
 - Pick a state at random and apply the appropriate backup
- Still need lots of computation, but does not get locked into hopelessly long sweeps
- Can you select states to backup intelligently? YES: an agent's experience can act as a guide.

Efficiency of DP

- ☐ To find an optimal policy is polynomial in the number of states...
- BUT, the number of states is often astronomical, e.g., often growing exponentially with the number of state variables (what Bellman called "the curse of dimensionality").
- □ In practice, classical DP can be applied to problems with a few millions of states.
- Asynchronous DP can be applied to larger problems, and is appropriate for parallel computation.
- □ It is surprisingly easy to come up with MDPs for which DP methods are not practical.

Summary

- Policy evaluation: backups without a max
- Policy improvement: form a greedy policy, if only locally
- Policy iteration: alternate the above two processes
- □ Value iteration: backups with a max
- □ Full backups (to be contrasted later with sample backups)
- Generalized Policy Iteration (GPI)
- □ Asynchronous DP: a way to avoid exhaustive sweeps
- **Bootstrapping**: updating estimates based on other estimates
- Biggest limitation of DP is that it requires a *probability model* (as opposed to a generative or simulation model)

Dynamic Programming Policy Evaluation



From Planning to Learning

- DP requires a *probability model* (as opposed to a generative or simulation model)
- We can interact with the world, learning a model (rewards and transitions) and then do DP
- **This approach is called model-based RL**
- **T** Full probability model may hard to learn though
- **Direct** learning of the value function from interaction
- □ Still focusing on evaluating a fixed policy

What if we don't have the transition probs?

Sample!

$$\square \mathbb{E}[f(S_{t+1}, r_{t+1}) | S_t] = \sum_{S_{t+1}} p(S_{t+1}, r_{t+1} | S_t) f(S_{t+1}, r_{t+1}) \approx \frac{1}{N} \sum_{k=1}^N f(S_{t_k+1}, r_{t_k+1})$$

Monte Carlo Methods

Monte Carlo is a neighbourhood of the country of Monaco
 Monte Carlo methods are named after the Monte Carlo Casino



Simple Monte Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha \Big[G_t - V(S_t) \Big]$$



Monte Carlo Methods

☐ Monte Carlo methods are learning methods Experience → values, policy

Monte Carlo methods can be used in two ways:

- *model-free:* No model necessary and still attains optimality
- *simulated:* Needs only a simulation, not a *full* model
- □ Monte Carlo methods learn from *complete* sample returns
 - Defined for episodic tasks (in the book)
- Like an associative version of a bandit method

Backup diagram for Monte Carlo

- **D** Entire rest of episode included
- Only one choice considered at each state (unlike DP)
 - thus, there will be an explore/exploit dilemma
- Does not bootstrap from successor states's values (unlike DP)
- Time required to estimate one state does not depend on the total number of states



Monte Carlo Policy Evaluation

- **Goal:** learn $v_{\pi}(s)$
- **Given:** some number of episodes under π which contain *s*
- **Idea:** Average returns observed after visits to s



- Every-Visit MC: average returns for every time s is visited in an episode
- □ *First-visit MC:* average returns only for *first* time *s* is visited in an episode
- **D** Both converge asymptotically

First-visit Monte Carlo policy evaluation

Initialize:

 $\pi \leftarrow \text{policy to be evaluated}$ $V \leftarrow \text{an arbitrary state-value function}$ $Returns(s) \leftarrow \text{an empty list, for all } s \in S$

Repeat forever:

Generate an episode using π For each state *s* appearing in the episode: $G \leftarrow$ return following the first occurrence of *s* Append *G* to Returns(s) $V(s) \leftarrow$ average(Returns(s))

MC vs supervised regression

- Target returns can be viewed as a supervised label (true value we want to fit)
- **I** State is the input
- ☐ We can use any function approximator to fit a function from states to returns! Neural nets, linear, nonparametric...
- Unlike supervised learning: there is strong correlation between inputs and between outputs!
- Due to the lack of iid assumptions, theoretical results from supervised learning cannot be directly applied

Blackjack example

Object: Have your card sum be greater than the dealer's without exceeding 21.

- States (200 of them):
 - current sum (12-21)
 - dealer's showing card (ace-10)
 - do I have a useable ace?



- **Reward:** +1 for winning, 0 for a draw, -1 for losing
- Actions: stick (stop receiving cards), hit (receive another card)
- **Policy:** Stick if my sum is 20 or 21, else hit
- **\square** No discounting ($\gamma = 1$)

Learned blackjack state-value functions



Simplest TD Method

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



TD methods bootstrap and sample

Bootstrapping: update involves an *estimate*

- MC does not bootstrap
- DP bootstraps
- TD bootstraps
- Sampling: update does not involve an expected value
 - MC samples
 - DP does not sample
 - TD samples

TD Prediction

Policy Evaluation (the prediction problem): for a given policy π , compute the state-value function v_{π}

Recall: Simple every-visit Monte Carlo method:

$$V(S_t) \leftarrow V(S_t) + \alpha \Big[G_t - V(S_t) \Big]$$

target: the actual return after time t

The simplest temporal-difference method TD(0):

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

target: an estimate of the return

Example: Driving Home

	Elapsed Time	Predicted	Predicted
State	(minutes)	Time to Go	Total Time
leaving office, friday at 6	0	30	30
reach car, raining	5	35	40
exiting highway	20	15	35
2ndary road, behind truck	30	10	40
entering home street	40	3	43
arrive home	43	0	43

Driving Home

Changes recommended by Monte Carlo methods (α =1) Changes recommended by TD methods (α =1)



Advantages of TD Learning

- TD methods do not require a model of the environment, only experience
- TD, but not MC, methods can be fully incremental
 - You can learn before knowing the final outcome
 - Less memory
 - Less peak computation
 - You can learn without the final outcome

From incomplete sequences

Both MC and TD converge (under certain assumptions to be detailed later), but which is faster?

Random Walk Example


TD and MC on the Random Walk



100 sequences of episodes

Batch Updating in TD and MC methods

Batch Updating: train completely on a finite amount of data, e.g., train repeatedly on 10 episodes until convergence.

Compute updates according to TD or MC, but only update estimates after each complete pass through the data.

For any finite Markov prediction task, under batch updating, TD converges for sufficiently small α .

Constant- α MC also converges under these conditions, but to a different answer!

Random Walk under Batch Updating



After each new episode, all previous episodes were treated as a batch, and algorithm was trained until convergence. All repeated 100 times.

You are the Predictor

Suppose you observe the following 8 episodes:



Assume Markov states, no discounting ($\gamma = 1$)

You are the Predictor





You are the Predictor

- The prediction that best matches the training data is V(A)=0
 - This minimizes the mean-square-error on the training set
 - This is what a batch Monte Carlo method gets
- If we consider the sequentiality of the problem, then we would set V(A)=.75
 - This is correct for the maximum likelihood estimate of a Markov model generating the data
 - i.e, if we do a best fit Markov model, and assume it is exactly correct, and then compute what it predicts (how?)
 - This is called the certainty-equivalence estimate
 - This is what TD gets

Application of TD Dopamine neuron activity modelling



Cf. Shultz, Dayan et al, 1996; and lots of follow-up work including MNI, Psych.

Summary so far

- Introduced one-step tabular model-free TD methods
- These methods bootstrap and sample, combining aspects of DP and MC methods
- TD methods are *computationally congenial*
- If the world is truly Markov, then TD methods will learn faster than MC methods
- MC methods have lower error on past data, but higher error on future data

Unified View



n-step TD Prediction



Mathematics of *n*-step TD Returns/Targets

• Monte Carlo: $G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T$

- **TD**: $G_t^{(1)} \doteq R_{t+1} + \gamma V_t(S_{t+1})$ • Use V_t to estimate remaining return
- *n*-step TD: • 2 step return: $G_t^{(2)} \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 V_t(S_{t+2})$

• *n*-step return: $G_t^{(n)} \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V_t(S_{t+n})$ with $[G_t^{(n)} \doteq G_t \text{ if } t+n \ge T]$

Forward View

Look forward from each state to determine update from future states and rewards:





n-step TD

• Recall the *n*-step return:

 $G_t^{(n)} \doteq R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n}), \quad n \ge 1, 0 \le t < T - n$

- Of course, this is <u>not available</u> until time t+n
- The natural algorithm is thus to wait until then:

$$V_{t+n}(S_t) \doteq V_{t+n-1}(S_t) + \alpha \left[G_t^{(n)} - V_{t+n-1}(S_t) \right], \qquad 0 \le t < T$$

• This is called *n*-step TD

n-step TD for estimating $V \approx v_{\pi}$

```
Initialize V(s) arbitrarily, s \in S
Parameters: step size \alpha \in (0, 1], a positive integer n
All store and access operations (for S_t and R_t) can take their index mod n
```

```
Repeat (for each episode):
    Initialize and store S_0 \neq terminal
   T \leftarrow \infty
   For t = 0, 1, 2, \ldots:
        If t < T, then:
            Take an action according to \pi(\cdot|S_t)
            Observe and store the next reward as R_{t+1} and the next state as S_{t+1}
            If S_{t+1} is terminal, then T \leftarrow t+1
       \tau \leftarrow t - n + 1 (\tau is the time whose state's estimate is being updated)
       If \tau > 0:
            G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n,T)} \gamma^{i-\tau-1} R_i
                                                                                                  G_{\tau}^{(n)}
            If \tau + n < T, then: G \leftarrow G + \gamma^n V(S_{\tau+n})
            V(S_{\tau}) \leftarrow V(S_{\tau}) + \alpha \left[ G - V(S_{\tau}) \right]
    Until \tau = T - 1
```



• How does 2-step TD work here?

• How about 3-step TD?

A Larger Example – 19-state Random Walk



- An intermediate α is best
- An intermediate *n* is best
- Do you think there is an optimal *n*? for every task?

Conclusions Regarding *n***-step Methods** (so far)

 Generalize Temporal-Difference and Monte Carlo learning methods, sliding from one to the other as *n* increases

• n = 1 is TD(0) $n = \infty$ is MC

- an intermediate *n* is often much better than either extreme
- applicable to both continuing and episodic problems
- There is some cost in computation
 - need to remember the last *n* states
 - learning is delayed by n steps
 - per-step computation is small and uniform, like TD

How to do control? GPI!

Generalized Policy Iteration (GPI):

any interaction of policy evaluation and policy improvement, independent of their granularity.



Monte Carlo Estimation of Action Values

Estimate q_{π} for the current policy π

$$\cdots \underbrace{S_t}_{S_t,A_t} \underbrace{R_{t+1}}_{S_{t+1}} \underbrace{S_{t+1}}_{S_{t+1},A_{t+1}} \underbrace{S_{t+2}}_{S_{t+2}} \underbrace{R_{t+3}}_{S_{t+3}} \underbrace{S_{t+3}}_{S_{t+3},A_{t+3}} \cdots$$

$$Q(S_t,A_t) \leftarrow Q(S_t,A_t) + \alpha(G_t - Q(S_t,A_t))$$

$$\text{where } G_t = \sum_{k=1}^{T-t} \gamma^{k-1} R_{t+k}$$

and *T* is the time of entering terminal state

Monte Carlo Estimation of Action Values (Q)

- \Box $q_{\pi}(s,a)$ average return starting from state *s* and action *a* following π
- Converges asymptotically *if* every state-action pair is visited
- Exploring starts: Every state-action pair has a non-zero probability of being the starting pair

On-policy Monte Carlo Control

On-policy: learn about policy currently executing
How do we get rid of exploring starts?

- The policy must be eternally *soft*:
 - $-\pi(a|s) > 0$ for all *s* and *a*
- e.g. ε-soft policy:

- probability of an action = $\frac{\epsilon}{|\mathcal{A}(s)|}$ or $1 - \epsilon + \frac{\epsilon}{|\mathcal{A}(s)|}$ non-max max (greedy)

- Similar to GPI: move policy *towards* greedy policy (e.g., ε-greedy)
- \square Converges to best ϵ -soft policy

On-policy MC Control

Initialize, for all $s \in S$, $a \in \mathcal{A}(s)$: $Q(s, a) \leftarrow \text{arbitrary}$ $Returns(s, a) \leftarrow \text{empty list}$ $\pi(a|s) \leftarrow \text{an arbitrary } \varepsilon\text{-soft policy}$

Repeat forever:

(a) Generate an episode using π (b) For each pair s, a appearing in the episode: $G \leftarrow$ return following the first occurrence of s, aAppend G to Returns(s, a) $Q(s, a) \leftarrow$ average(Returns(s, a)) (c) For each s in the episode: $A^* \leftarrow$ arg max_a Q(s, a)For all $a \in \mathcal{A}(s)$: $\pi(a|s) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon/|\mathcal{A}(s)| & \text{if } a = A^* \\ \varepsilon/|\mathcal{A}(s)| & \text{if } a \neq A^* \end{cases}$

TD-Style Learning for Action-Values

Estimate q_{π} for the current policy π

$$\cdots \underbrace{S_{t}}_{S_{t},A_{t}} \underbrace{R_{t+1}}_{S_{t+1}} \underbrace{S_{t+1}}_{S_{t+1},A_{t+1}} \underbrace{R_{t+2}}_{S_{t+2}} \underbrace{R_{t+3}}_{S_{t+2},A_{t+2}} \underbrace{S_{t+3}}_{S_{t+3},A_{t+3}} \cdots$$

After every transition from a nonterminal state, S_t , do this: $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \Big[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t) \Big]$ If S_{t+1} is terminal, then define $Q(S_{t+1}, A_{t+1}) = 0$ Turn this into a control method by always updating the policy to be greedy with respect to the current estimate:

Windy Gridworld



undiscounted, episodic, reward = -1 until goal

Results of Sarsa on the Windy Gridworld



Q-Learning: Off-Policy TD Control

One-step Q-learning:

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

 $\begin{array}{ll} \mbox{Initialize } Q(s,a), \forall s \in \mathbb{S}, a \in \mathcal{A}(s), \mbox{ arbitrarily, and } Q(terminal-state, \cdot) = 0 \\ \mbox{Repeat (for each episode):} \\ \mbox{Initialize } S \\ \mbox{Repeat (for each step of episode):} \\ \mbox{Choose } A \mbox{ from } S \mbox{ using policy derived from } Q \mbox{ (e.g., ε-greedy)} \\ \mbox{Take action } A, \mbox{ observe } R, S' \\ Q(S,A) \leftarrow Q(S,A) + \alpha [R + \gamma \max_a Q(S',a) - Q(S,A)] \\ S \leftarrow S'; \\ \mbox{until } S \mbox{ is terminal} \end{array}$

Cliffwalking



Expected Sarsa

Instead of the sample value-of-next-state, use the expectation!

$$Q(S_{t}, A_{t}) \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \mathbb{E}[Q(S_{t+1}, A_{t+1}) \mid S_{t+1}] - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$



• Expected Sarsa's performs better than Sarsa (but costs more)

Performance on the Cliff-walking Task



R. S. Sutton and A. G. Barto: Reinforcement Learning: An Introduction

Off-policy Expected Sarsa

- Expected Sarsa generalizes to arbitrary behavior policies μ
 - in which case it includes Q-learning as the special case in which π is the greedy policy

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

$$\begin{array}{c} \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big] \\ \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$

This idea seems to be new

Maximization Bias Example



Tabular Q-learning: $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \Big[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \Big]$

Hado van Hasselt 2010

Double Q-Learning

- Train 2 action-value functions, Q_1 and Q_2
- Do Q-learning on both, but
 - never on the same time steps (Q_1 and Q_2 are indep.)
 - pick Q_1 or Q_2 at random to be updated on each step
- If updating Q_1 , use Q_2 for the value of the next state: $Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha \Big(R_{t+1} + Q_2 \big(S_{t+1}, \operatorname*{arg\,max}_a Q_1(S_{t+1}, a) \big) - Q_1(S_t, A_t) \Big)$
- Action selections are (say) ε -greedy with respect to the sum of Q_1 and Q_2

Hado van Hasselt 2010

Double Q-Learning

 $\begin{array}{l} \mbox{Initialize } Q_1(s,a) \mbox{ and } Q_2(s,a), \forall s \in \mathbb{S}, a \in \mathcal{A}(s), \mbox{ arbitrarily} \\ \mbox{Initialize } Q_1(terminal-state, \cdot) = Q_2(terminal-state, \cdot) = 0 \\ \mbox{Repeat (for each episode):} \\ \mbox{Initialize } S \\ \mbox{Repeat (for each step of episode):} \\ \mbox{Choose } A \mbox{ from } S \mbox{ using policy derived from } Q_1 \mbox{ and } Q_2 \mbox{ (e.g., } \varepsilon\mbox{-greedy in } Q_1 + Q_2) \\ \mbox{Take action } A, \mbox{ observe } R, S' \\ \mbox{With } 0.5 \mbox{ probabilility:} \\ \mbox{} Q_1(S,A) \leftarrow Q_1(S,A) + \alpha \Big(R + \gamma Q_2 \big(S', \mbox{arg max}_a Q_1(S',a) \big) - Q_1(S,A) \Big) \\ \mbox{else:} \\ \mbox{} Q_2(S,A) \leftarrow Q_2(S,A) + \alpha \Big(R + \gamma Q_1 \big(S', \mbox{arg max}_a Q_2(S',a) \big) - Q_2(S,A) \Big) \\ \mbox{} S \leftarrow S'; \\ \mbox{ until } S \mbox{ is terminal} \end{array} \right)$

Example of Maximization Bias



Double Q-learning:

$$Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha \Big[R_{t+1} + \gamma Q_2 \big(S_{t+1}, \operatorname*{arg\,max}_a Q_1(S_{t+1}, a) \big) - Q_1(S_t, A_t) \Big]$$

Summary

- Introduced one-step tabular model-free TD methods
- These methods bootstrap and sample, combining aspects of DP and MC methods
- TD methods are *computationally congenial*
- If the world is truly Markov, then TD methods will learn faster than MC methods
- MC methods have lower error on past data, but higher error on future data
- Extend prediction to control by employing some form of GPI
 - On-policy control: Sarsa, Expected Sarsa
 - Off-policy control: Q-learning, Expected Sarsa
- Avoiding maximization bias with Double Q-learning