Reinforcement learning Master CPS, Year 2 Semester 1

Lucian Buşoniu, Florin Gogianu



Part III

Temporal differences

Exact reinforcement learning



Reinforcement learning problem

Temporal differences

- Optimal solution
- Exact dynamic programming
- Exact reinforcement learning
- Approximation techniques
- Approximate dynamic programming
- Approximate reinforcement learning



Algorithm landscape

Monte Carlo

By model usage:

- Model-based: f, ρ known a priori
- Model-free: f, ρ unknown (reinforcement learning)

By interaction level:

- Offline: algorithm runs in advance
- Online: algorithm runs with the system

Exact vs. approximate:

- Exact: x, u small number of discrete values
- Approximate: x, u continuous (or many discrete values)



Supervised learning

Monte Carlo

Reinforcement learning

Temporal differences

Unsupervised learning

more informative feedback

less informative feedback

- Supervised: for each training sample, correct output known
- Unsupervised: only input samples, no outputs; find patterns in the data
- Reinforcement: correct actions not available, only rewards

But note: RL finds dynamical optimal control!



- Monte Carlo, MC
- **Exploration**
- Temporal differences, TD
- Accelerating TD methods
- Recap



- Monte Carlo, MC
- Exploration

- Accelerating TD methods
- Recap



Reminder: Policy iteration

Policy iteration with Q-functions

```
initialize policy h_0 arbitrarily repeat at each iteration \ell
1: policy evaluation: find Q^{h_\ell}
2: policy improvement:
find h_{\ell+1}(x) = \arg\max_u Q^{h_\ell}(x,u)
until convergence to h^*
```

Temporal differences

Note: In RL, we generally use **Q-functions** so policy improvement does not require a model



Policy evaluation

Monte Carlo

000000

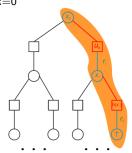
To find Q^h :

- So far: model-based methods
- Reinforcement learning: model not available
- Learn Q^h from offline data or via online interaction with the system



"Monte-Carlo" policy evaluation

Recall:
$$Q^h(x_0, u_0) = \sum_{k=0}^{\infty} \gamma^k r_{k+1}$$

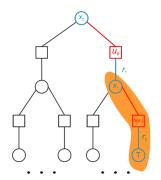


- Trajectory from (x_0, u_0) to x_K (terminal) using $u_1 = h(x_1)$, $u_2 = h(x_2)$, etc.
- $\Rightarrow Q^h(x_0, u_0) = \text{return along trajectory:}$

$$Q^h(x_0, u_0) = \sum_{i=0}^{K-1} \gamma^j r_{j+1}$$



0000000

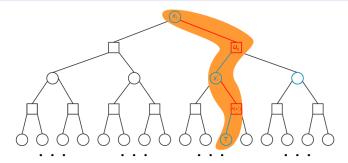


• Moreover, at each step k:

$$Q^{h}(x_{k}, u_{k}) = \sum_{j=k}^{K-1} \gamma^{j-k} r_{j+1}$$



Monte-Carlo policy evaluation: Stochastic case



- N trajectories (differing due to stochastic transitions)
- Estimated Q value = **mean** of the returns, e.g.

$$Q^{h}(x_{0}, u_{0}) = \frac{1}{N} \sum_{i=1}^{N} \sum_{j=0}^{K_{i}-1} \gamma^{j} r_{i,j+1}$$



end for

Monte-Carlo policy iteration

```
Monte-Carlo policy iteration
   for each iteration ℓ do
        perform N trajectories applying h_{\ell}
        reset to 0 accumulator A(x, u), counter C(x, u)
       for each step k of each trajectory i do
            A(x_k, u_k) \leftarrow A(x_k, u_k) + \sum_{i=k}^{K_i-1} \gamma^{j-k} r_{i,i+1} (return)
            C(x_k, u_k) \leftarrow C(x_k, u_k) + 1
       end for
       Q^{h_{\ell}}(x,u) \leftarrow A(x,u)/C(x,u)
        h_{\ell+1}(x) \leftarrow \operatorname{arg\,max}_{\ell} Q^{h_{\ell}}(x, u)
```

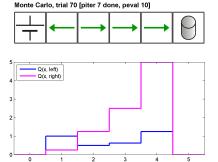
Note: must guarantee that terminal state is reached!



000000

2

0 -



state, x

policy iteration

—● h - h*



Accelerating TD

- Monte Carlo, MC
- 2 Exploration
- Accelerating TD methods
- Recap



The need for exploration

In the MC estimate:

$$Q^h(x, u) \leftarrow A(x, u)/C(x, u)$$

how to ensure C(x, u) > 0 – information about each (x, u)?

Temporal differences

- Initial states x₀ representative
- Actions:

 u_0 representative, sometimes different from $h(x_0)$ and additionally, possibly: u_k representative, sometimes different from $h(x_k)$



Exploration-exploitation dilemma

- Exploration is necessary: actions different from the current policy
- Exploitation of current knowledge is necessary: current policy must be applied for good performance

Temporal differences

The exploration-exploitation dilemma

- essential in all RL algorithms

(not only in MC)



ε -greedy strategy

 Simple solution to the exploration-exploitation dilemma: ε -greedy

Temporal differences

$$u_k = \begin{cases} h(x_k) = \arg\max_u Q(x_k, u) & \text{with probability } (1 - \varepsilon_k) \\ \text{a uniformly random action} & \text{w.p. } \varepsilon_k \end{cases}$$

- Exploration probability $\varepsilon_k \in (0,1)$ usually decreased over time
- Main disadvantage: when exploring, actions are fully random, leading to poor performance



Softmax strategy

Monte Carlo

Action selection:

$$u_k = u$$
 w.p. $\frac{e^{Q(x_k,u)/ au_k}}{\sum_{u'} e^{Q(x_k,u')/ au_k}}$

where $\tau_k > 0$ is the exploration temperature

Temporal differences

- Taking $\tau \to 0$, greedy selection recovered; $\tau \to \infty$ gives uniform random
- Compared to ε -greedy, better actions are more likely to be applied even when exploring



Bandit-based exploration

Monte Carlo



Temporal differences

At single state, exploration modeled as **multi-armed bandit**:

- Action $j = \text{arm with reward distribution } \rho_i$, expectation μ_i
- Best arm (optimal action) has expected value μ^*
- At step k, we pull arm (try action) j_k , getting $r_k \sim \rho_{i_k}$
- Objective: After *n* pulls, small regret: $\sum_{k=1}^{n} \mu^* \mu_{i_k}$



UCB algorithm

Monte Carlo

Often-used algorithm: after *n* steps, pick arm with largest upper confidence bound:

$$b(j) = \widehat{\mu}_j + \sqrt{c \frac{\log n}{n_j}}$$

where:

- $\hat{\mu}_i$ = mean of rewards observed for arm i so far
- n_i = how many times arm i was pulled
- c tunable constant, e.g. 3/2

These are only a few simple methods, many others exist, e.g. Bayesian exploration, intrinsic rewards, optimistic initialization etc.



Exploration

Monte Carlo

- Temporal differences, TD
 - Introduction
 - SARSA
 - Q-learning
 - Discussion
- Accelerating TD methods
- 5 Recap



Recap

Monte-Carlo with incremental updates

Consider the return sample from step *k* onwards:

$$R_k = \sum_{j=k}^{K-1} \gamma^{j-k} r_{j+1}$$

Instead of averaging such samples to get Q, perform incremental updates:

$$Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k [R_k - Q(x_k, u_k)]$$

where α_k is a step size, or learning rate



Discussion

 Incremental MC motivated by time-varying problems (recent samples have larger weights), but works in time-invariant MDPs as well

Temporal differences

- No longer need to store accumulators A and counters C, just directly the Q-values
- If α satisfies "stochastic-approximation" conditions:
 - **1** decreases to 0, $\sum_{k=0}^{\infty} \alpha_k^2 = \text{finite}$
 - 2 but not too quickly, $\sum_{k=0}^{\infty} \alpha_k \to \infty$

method converges: lim when # of samples $\rightarrow \infty$ = Q-value

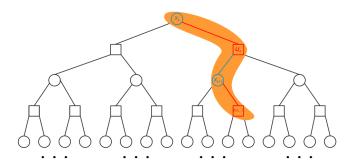


From Monte-Carlo to temporal differences

To avoid waiting until trajectory finishes, recall Bellman equation:

$$Q^h(x,u)=\mathrm{E}_{x'}\left\{\tilde{\rho}(x,u,x')+\gamma Q^h(x',h(x'))\right\}$$
 and use the return estimate:

$$\hat{R}_k = r_{k+1} + \gamma Q(x_{k+1}, h(x_{k+1}))$$





Otherwise, update remains the same, but let us make the return estimate explicit:

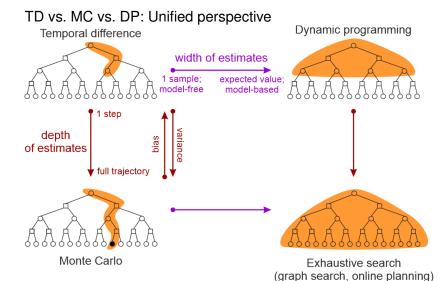
Temporal differences

$$Q(x_{k}, u_{k}) \leftarrow Q(x_{k}, u_{k}) + \alpha_{k} [\hat{R}_{k} - Q(x_{k}, u_{k})]$$

= $Q(x_{k}, u_{k}) + \alpha_{k} [r_{k+1} + \gamma Q(x_{k+1}, h(x_{k+1})) - Q(x_{k}, u_{k})]$

- [...] is the **temporal difference** between two estimates of $Q(x_k, u_k)$, using information at subsequent time steps
- Model-free, data-based updates (like MC): r_{k+1} , x_{k+1} e.g. observed while interacting online
- Updates estimate $Q(x_k, u_k)$ using another estimate, $Q(x_{k+1}, h(x_{k+1}))$: **bootstrapping**
- Dynamic programming also bootstraps, but using a model







TD for policy evaluation

```
for each trajectory do
    initialize x_0, choose the initial action u_0
    repeat at each step k
        apply u_k, measure x_{k+1}, receive r_{k+1}
        choose the next action u_{k+1} \sim h(x_{k+1})
         Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k
                    [r_{k+1} + \gamma Q(x_{k+1}, u_{k+1}) - Q(x_k, u_k)]
    until trajectory finished
end for
```

Note: we replaced $h(x_{k+1})$ by u_{k+1} , chosen according to h



choose the next action $u_{k+1} \sim h(x_{k+1})$

• Information about $(x, u) \neq (x, h(x))$ necessary \Rightarrow exploration

Temporal differences

- h must be followed \Rightarrow exploitation
- E.g. ε -greedy:

$$u_{k+1} = \begin{cases} h(x_{k+1}) & \text{w.p. } (1 - \varepsilon_{k+1}) \\ \text{unif. random} & \text{w.p. } \varepsilon_{k+1} \end{cases}$$



2 Exploration

Monte Carlo

- Temporal differences, TD
 - Introduction
 - SARSA
 - Q-learning
 - Discussion
- 4 Accelerating TD methods
- 5 Recap



Recall: MC

Monte Carlo

MC policy iteration

```
for each iteration \ell do
     perform N trajectories applying h_{\ell}
     reset to 0 accumulator A(x, u), counter C(x, u)
    for each step k of each trajectory i do
         A(x_k, u_k) \leftarrow A(x_k, u_k) + \sum_{i=k}^{K_i-1} \gamma^{j-k} r_{i,i+1} (return)
         C(x_k, u_k) \leftarrow C(x_k, u_k) + 1
    end for
     Q^{h_{\ell}}(x,u) \leftarrow A(x,u)/C(x,u)
     h_{\ell+1}(x) \leftarrow \operatorname{arg\,max}_{\ell} Q^{h_{\ell}}(x, u)
end for
```



Optimistic policy improvement

- Policy unchanged for N trajectories
- ⇒ Algorithm learns slowly
 - Policy improvement after each trajectoryoptimistic
 - We will also use ε -greedy exploration



Optimistic MC

Optimistic MC

```
initialize to 0 accumulator A(x, u), counter C(x, u)
for each trajectory do
     execute trajectory, e.g., applying \varepsilon-greedy:
     u_k = \begin{cases} \arg \max_u Q(x_k, u) & \text{w.p. } (1 - \varepsilon_k) \\ \text{unif. random} & \text{w.p. } \varepsilon_k \end{cases}
     for each step k do
          A(x_k, u_k) \leftarrow A(x_k, u_k) + \sum_{j=k}^{K-1} \gamma^{j-k} r_{j+1}
           C(x_k, u_k) \leftarrow C(x_k, u_k) + 1
     end for
     Q(x, u) \leftarrow A(x, u)/C(x, u)
end for
```

Temporal differences

- h implicit, greedy in Q
- update of $Q \Rightarrow$ implicit improvement of h



Optimism in TD

Monte Carlo

- Earlier TD algorithm: fixed h
- What is the fastest we can improve h in TD? After each transition.

Temporal differences

- ⇒ interpretation: policy iteration optimistic at the transition level
 - h implicit, greedy in Q (updating $Q \Rightarrow$ implicit improvement of h)



SARSA

Monte Carlo

```
SARSA with \varepsilon-greedy
    for each trajectory do
          initialize x₀
          u_0 = \begin{cases} \arg \max_u Q(x_0, u) & \text{w.p. } (1 - \varepsilon_0) \\ \text{unif. random} & \text{w.p. } \varepsilon_0 \end{cases}
          repeat at each step k
                 apply u_k, measure x_{k+1}, receive r_{k+1}
                u_{k+1} = \begin{cases} \arg\max_{u} Q(x_{k+1}, u) & \text{w.p. } (1 - \varepsilon_{k+1}) \\ \text{unif. random} & \text{w.p. } \varepsilon_{k+1} \end{cases}
                  Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k
                                   [r_{k+1} + \gamma Q(x_{k+1}, u_{k+1}) - Q(x_k, u_k)]
          until trajectory finished
    end for
```



Origin of name "SARSA"

Monte Carlo

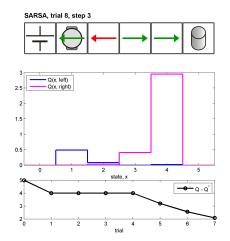
$$(x_k, u_k, r_{k+1}, x_{k+1}, u_{k+1}) =$$

(State, Action, Reward, State, Action) = SARSA

Temporal differences



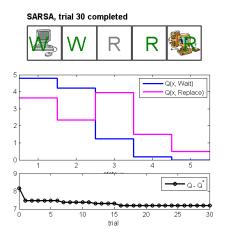
Parameters: $\alpha = 0.2$, $\varepsilon = 0.3$ (constant) $x_0 = 2$ or 3 (random)





Parameters: $\alpha = 0.1$, $\varepsilon = 0.3$ (constant), 20 steps per trajectory $x_0 = 1$

Temporal differences





- Monte Carlo, MC

- Temporal differences, TD
 - Introduction
 - SARSA
 - Q-learning
- Accelerating TD methods



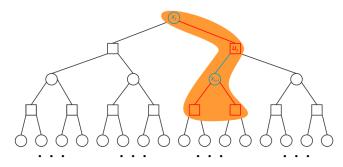
Bootstrapping estimate of Q^*

Bellman optimality equation:

$$Q^*(x,u) = \mathrm{E}_{x'} \left\{ \tilde{\rho}(x,u,x') + \gamma \max_{u'} Q^*(x',u') \right\}$$

leads to estimate:

$$\hat{Q}_k = r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u')$$





$$Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k [\hat{Q}_k - Q(x_k, u_k)]$$
$$[r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u') - Q(x_k, u_k)]$$



Q-learning with ε -greedy

Q-learning

Monte Carlo

for each trajectory do initialize x_0 repeat at each step k $u_k = \begin{cases} \arg\max_u Q(x_k, u) & \text{w.p. } (1 - \varepsilon_k) \\ \text{unif. random} & \text{w.p. } \varepsilon_k \\ \text{apply } u_k, \text{ measure } x_{k+1}, \text{ receive } r_{k+1} \end{cases}$

 $[r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u') - Q(x_k, u_k)]$

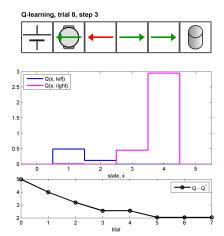
until trajectory finished **end for**

 $Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k$



Cleaning robot: Q-learning, demo

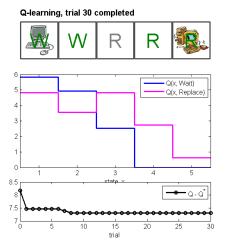
Parameters – same as SARSA: $\alpha =$ 0.2, $\varepsilon =$ 0.3 (constant) $x_0 =$ 2 or 3 (random)





Machine replacement: Q-learning, demo

Parameters: $\alpha = 0.1$, $\varepsilon = 0.3$ (constant), 20 steps per trajectory $x_0 = 1$





- Monte Carlo, MC
- 2 Exploration
- Temporal differences, TD
 - Introduction
 - SARSA
 - Q-learning
 - Discussion
- 4 Accelerating TD methods
- 5 Recap



Convergence

Monte Carlo

Conditions for convergence to *Q**:

- \bullet All pairs (x, u) continue to be updated: ensured by **exploration**, e.g. ε -greedy
- Stochastic-approximation conditions: learning rate decreases to zero, $\sum_{k=0}^{\infty} \alpha_k^2 = \text{finite}$, but not too quickly: $\sum_{k=0}^{\infty} \alpha_k \to \infty$

Additionally, for SARSA:

The policy must become greedy at infinity e.g. $\lim_{k\to\infty} \varepsilon_k = 0$



On-policy / off-policy

SARSA: on-policy

Always estimates the Q-function of the current policy

Temporal differences

Q-learning: off-policy

 Independently of the current policy, always estimates the optimal Q-function



TD: Discussion

Advantages

- Simple to understand and implement
- Low complexity ⇒ fast execution

SARSA vs. Q-learning

 SARSA less complex than Q-learning (no max in the Q-function update)

Learning rate and exploration sequences α_k , ε_k significantly influence performance

Main disadvantage

Large amount of data required



- Monte Carlo, MC
- Exploration
- Temporal differences, TD
- Accelerating TD methods
 - Motivation
 - Experience replay
 - n-step returns
- 6 Recap



The need to accelerate TD

Main disadvantage: TD learns slowly - requires a lot of data

In practice, data costs:

time

Monte Carlo

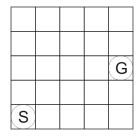
- profit (low performance due to exploration)
- system wear

Accelerating RL = efficient use of data



Example: 2D navigation

Monte Carlo

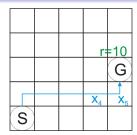


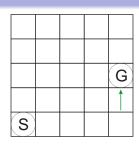
- Navigation in a discrete 2D world from Start to Goal
- Reward = 10 only upon reaching G (terminal state)



Example: TD

Monte Carlo





- We choose SARSA, $\alpha = 1$; initialize Q = 0
- Updates along trajectory on the left:

. . .

$$Q(x_4, u_4) = 0 + \gamma \cdot Q(x_5, u_5) = 0$$

 $Q(x_5, u_5) = 10 + \gamma \cdot 0 = 10$

A new transition from x₄ to x₅
 needed to propagate the information to x₄!



Accelerating TD: 2 ideas

- Store and replay experience
- Use n-step returns



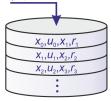
- Monte Carlo, MC
- 2 Exploration
- Temporal differences, TD
- Accelerating TD methods
 - Motivation
 - Experience replay
 - n-step returns
- 5 Recap



Experience replay (ER)

Monte Carlo

• Store each transition $(x_k, u_k, x_{k+1}, r_{k+1})$ (and for SARSA, also u_{k+1}) in a database



• At each step, **replay** *m* transitions from database (in addition to normal updates)



Q-learning with ER

Monte Carlo

Q-learning with ER

```
for each trajectory do
    initialize x₀
    repeat at each step k
        apply u_k, measure x_{k+1}, receive r_{k+1}
         Q(x_k, u_k) \leftarrow Q(x_k, u_k) + \alpha_k
                    [r_{k+1} + \gamma \max_{u'} Q(x_{k+1}, u') - Q(x_k, u_k)]
        add (x_k, u_k, x_{k+1}, r_{k+1}) to the database
        ReplayExperience
    until trajectory ends
end for
```



ReplayExperience procedure

ReplayExperience

Monte Carlo

```
loop m times
    fetch a transition (x, u, x', r) from the database
    Q(x, u) \leftarrow Q(x, u) + \alpha
               [r + \gamma \max_{u'} Q(x', u') - Q(x, u)]
end loop
```



Temporal differences

Monte Carlo

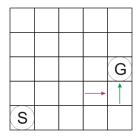
Replay direction

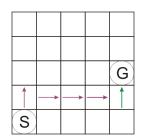
Order of replaying transitions:

- Forward
- Backward
- Arbitrary



Example: Influence of replay direction





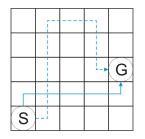
• Green: normal updates, purple: experience replay

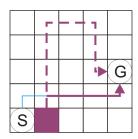
Temporal differences

- Left: forward replay; right: backward replay
- Backward replay preferable in exact RL



Example: Aggregating information



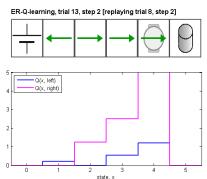


- Experience replay aggregates information from multiple trajectories
- The indicated cell benefits from information along both trajectories



Parameters: $\alpha = 0.2$, $\varepsilon = 0.3$, m = 5, backward direction $x_0 = 2 \text{ or } 3 \text{ (random)}$

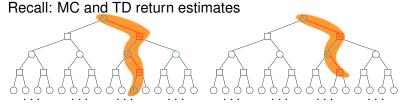
Temporal differences





- Monte Carlo, MC
- Exploration
- Temporal differences, TD
- Accelerating TD methods
 - Motivation
 - Experience replay
 - n-step returns
- 5 Recap





Temporal differences

$$R_{k} = \sum_{j=k}^{K-1} \gamma^{j-k} r_{j+1}$$

$$\hat{R}_{k} = r_{k+1} + \gamma Q(x_{k+1}, h(x_{k+1}))$$

Is there something in-between?

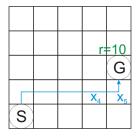


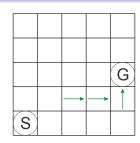


SARSA (on-policy):

$$\hat{R}_{k} = r_{k+1} + \gamma r_{k+2} + \dots + \gamma^{n-1} r_{k+n} + \gamma^{n} Q(x_{k+n}, u_{k+n})$$







For n = 3:

Monte Carlo

$$Q(x_5,u_5) = 10+0$$
 (terminal) $Q(x_4,u_4) = 0+\gamma 10+0$ (terminal) $Q(x_3,u_3) = 0+\gamma 0+\gamma^2 10+0$ (terminal) $Q(x_2,u_2) = 0+\gamma 0+\gamma^2 0+\gamma^3 0$ (bootstrap)

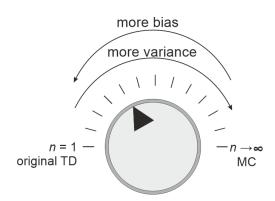
Temporal differences

. . .



- n = 1 recovers TD, $n \to \infty$ recovers MC
- Intermediate values mix TD and MC, leading to a tunable bias-variance tradeoff

Temporal differences





- Monte Carlo, MC
- Exploration

- Accelerating TD methods
- Recap



Recap: Methods in Part III

Monte-Carlo methods, MC:

- MC policy iteration
- MC with incremental updates

Exploration-exploitation dilemma:

- ε -greedy widely used
- Many other solutions exist, like UCB

Temporal differences

Temporal differences, TD:

- TD for policy evaluation
- Optimistic policy improvements
- SARSA
- Q-learning

Accelerating TD:

- Experience replay
- n-step returns



Key terms in this part

Monte Carlo

- Monte-Carlo methods
- exploration-exploitation dilemma
- \bullet ϵ -greedy, softmax, bandits
- optimistic policy improvement
- learning rate
- bootstrapping
- temporal differences
- SARSA
- Q-learning
- experience replay
- n-step returns



Exercises

Monte Carlo

• Does the exponential schedule $\alpha_k = \alpha^k$, with $\alpha \in (0,1)$ a constant, satisfy the stochastic approximation conditions?

Temporal differences

- 2 Is Q-learning guaranteed to converge when $\varepsilon_k = \varepsilon$, a constant in (0, 1)? What about SARSA? How about when you use an exponential decrease $\varepsilon_k = \varepsilon^k$?
- Would a Monte-Carlo algorithm that improves the policy after every transition (like TD) make sense?
- Would Q-learning (without n-step returns as they are nontrivial in the off-policy case) propagate information faster than SARSA for the gridworld trajectory example?



Exercises (cont'd)

Monte Carlo

- Assuming that we have access to a model only for the purposes of policy improvement, provide V-function alternates for all algorithms in this part. Do this in the same order as for Q-functions:
 - Monte Carlo estimates, averaging-based and incremental
 - Bootstrapping estimates and updates
 - Policy evaluation, SARSA, and Q-learning

Don't forget to draw trees and highlights, it will help you visualize things.

