

Today's central concepts

Previously, we have assumed that we could derive a model as well as obtain numerical values for the parameters.

What if everything is unknown?

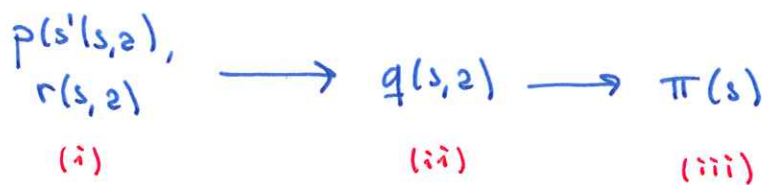
Learn by interacting with the system!

In Monte-Carlo methods:

(*):
(or several)

- i, a ^(*) trajectory is generated under some policy,
- ii, new estimates are computed and the policy is updated,
- iii, repeat to i.

We can focus to estimate different quantities:



(i): Directly compute $\hat{p}(s'|s,a)$ and $\hat{r}(s,a)$, then use standard DP to compute policy.

- Contains much information regarding the system
- Need to estimate $|S|^2 \times |A| + |S| \times |A|$ numbers.

(ii): The state-action function $q(s, a)$ is enough if we only want to control the system

- $|S| \times |A|$ numbers to estimate
- Idea: Keep track of rewards (to go) observed in (state, action)-pairs.

(iii): Actually, if our goal is to control the system, we could directly focus on $\pi(s)$. We'll do it in the next session (policy gradient).

In online/incremental/TD learning-methods,

the estimates as well as the policy are updated in each time-step.

(We don't wait for an episode to finish.)

Note:
Some problems don't even have episodes!

- Q-learning:

- Pick an action a (e.g. ϵ -greedy)
- Apply a in s and observe $r(s, a)$ and s'
- Update

$$q(s, a) \leftarrow q(s, a) + \alpha [r(s, a) + \lambda \max_{a'} q(s', a') - q(s, a)] =$$

reward-to-go from s using a
 $(1-\alpha) q(s, a)$ old estimate
 $+ \alpha [r(s, a) + \lambda \max_{a'} q(s', a')]$ new information, immediate reward
 $=$

$\underbrace{\hspace{10em}}$ assume greedy action in next time-step

iv, Repeat to i).

Under suitable assumptions.
See lecture.

Note:

- $q(s, a)$ will converge to $q^*(s, a)$.
- The actual action we implement in the next time-step is not necessarily the greedy!
- Off-policy algorithm since we learn value of π^* , but that's not the policy we apply to the system.

- SARSA:

- Take action a and observe $r(s, a)$ and s' .
- Pick an action a' (from s') (e.g. ϵ -greedy).
- Update

$$q(s, a) \leftarrow q(s, a) + \alpha [r(s, a) + \lambda \underbrace{q(s', a')}_{\text{(estimated) reward-to-go from the next time-step}} - q(s, a)]$$

iv, Set $a \leftarrow a'$ and repeat from i).

Note:

- SARSA learns the optimal policy taking into account how we actually select actions (that we explore).
- On-policy since we learn value of the policy we implement.

Ex. 4.3

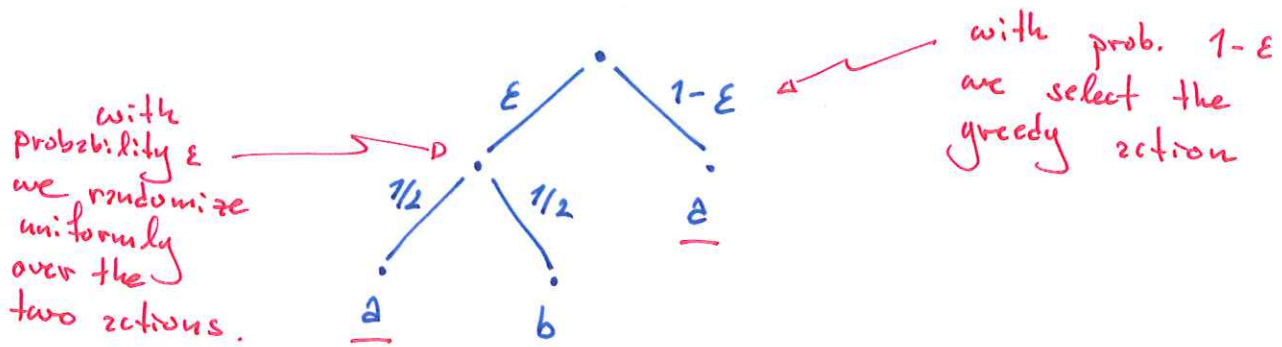
In ϵ -greedy action selection, what is the probability that the greedy action is selected if:

a, $\epsilon = 0.5$ and $|\mathcal{A}| = 2$?

b, $\epsilon = 1/5$ and $|\mathcal{A}| = 3$?

Solution:

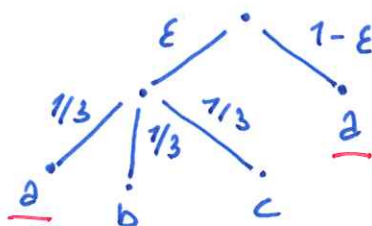
a, Assume $\mathcal{A} = \{a, b\}$ and a is the greedy action.



The greedy action can be reached with total probability:

$$\epsilon \cdot \frac{1}{2} + (1 - \epsilon) = \left\{ \epsilon = \frac{1}{2} \right\} = \frac{1}{4} + 1 - \frac{1}{2} = \frac{3}{4}.$$

b, Assume $\mathcal{A} = \{a, b, c\}$ and that a is greedy:



The total probability is now: $\epsilon \cdot \frac{1}{3} + (1 - \epsilon) = \left\{ \epsilon = \frac{1}{5} \right\} = \frac{1}{5} \cdot \frac{1}{3} + \frac{4}{5} = \frac{13}{5}.$

Consider the following observed trajectory:

time	current state	reward	action	next state
1	s_1	-2	a_1	s_3
2	s_3	6	a_1	s_3
3	s_3	4	a_2	s_2
4	s_2	-2	a_1	s_2
5	s_2	2	a_2	s_1

a, Perform Q-learning with $\gamma = 0.8$ and $\alpha = 0.5$.

b, what is the system's optimal policy, assuming that the algorithm has converged after these five steps? (note: it has not.)

Solution:

a, Recall the Q-learning update equation:

$$q(s, a) \leftarrow q(s, a) + \alpha [r(s, a) + \gamma \max_{a'} q(s', a') - q(s, a)]$$

Assume we initialize the q-table with zeroes.

Then:

	a_1	a_2
s_1	0	0
s_2	0	0
s_3	0	0

$q(s, a) = q(s, a_1) \leftarrow$

$$q(s, a) + \alpha [r(s, a) + \gamma \max_{a'} q(s', a') - q(s, a)] =$$

$$\underbrace{q(s, a_1)}_{=0} + \alpha [\underbrace{r(s, a_1)}_{=-2} + \gamma \underbrace{\max_{a'} q(s_3, a')}_{=0} - \underbrace{q(s, a_1)}_{=0}] =$$

for all $a' = 0$

$$0 + 0.5[-2 + 0 - 0] = -1$$

$s = s_1$
 $a = a_1$
 $r = -2$
 $s' = s_3$

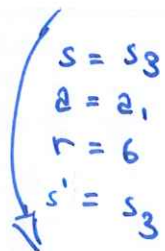
	a_1	a_2
s_1	-1	0
s_2	0	0
s_3	0	0

$q(s_3, a_1) \leftarrow$

$$q(s_3, a_1) + \alpha [r(s_3, a_1) + \lambda \max_{a'} q(s_3, a') - q(s_3, a_1)] =$$

$\underbrace{q(s_3, a_1)}_{=0} + \alpha \left[\underbrace{r(s_3, a_1)}_{=6} + \lambda \max_{a'} \underbrace{q(s_3, a')}_{=0 \text{ for all } a'} - \underbrace{q(s_3, a_1)}_{=0} \right] =$

$0.5 \cdot 6 = 3$



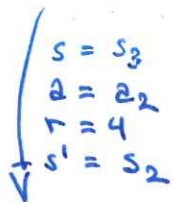
	a_1	a_2
s_1	-1	0
s_2	0	0
s_3	3	0

$q(s_3, a_2) \leftarrow$

$$q(s_3, a_2) + \alpha [r(s_3, a_2) + \lambda \max_{a'} q(s_2, a') - q(s_3, a_2)] =$$

$\underbrace{q(s_3, a_2)}_{=0} + \alpha \left[\underbrace{r(s_3, a_2)}_{=4} + \lambda \max_{a'} \underbrace{q(s_2, a')}_{=0 \text{ for all } a'} - \underbrace{q(s_3, a_2)}_{=0} \right] =$

$0.5 \times 4 = 2$



	a_1	a_2
s_1	-1	0
s_2	0	0
s_3	3	2

$q(s_2, a_1) \leftarrow$

$$q(s_2, a_1) + \alpha [r(s_2, a_1) + \lambda \max_{a'} q(s_2, a') - q(s_2, a_1)] =$$

$\underbrace{q(s_2, a_1)}_{=0} + \alpha \left[\underbrace{r(s_2, a_1)}_{=-2} + \lambda \max_{a'} \underbrace{q(s_2, a')}_{=0 \text{ for all } a'} - \underbrace{q(s_2, a_1)}_{=0} \right] =$

$0.5 \cdot (-2) = -1$



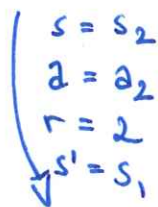
	a_1	a_2
s_1	-1	0
s_2	-1	0
s_3	3	2

$q(s_2, a_2) \leftarrow$

$$q(s_2, a_2) + \alpha [r(s_2, a_2) + \lambda \max_{a'} q(s_1, a') - q(s_2, a_2)] =$$

$\underbrace{q(s_2, a_2)}_{=0} + \alpha \left[\underbrace{r(s_2, a_2)}_{=2} + \lambda \max_{a'} \underbrace{q(s_1, a')}_{=0} - \underbrace{q(s_2, a_2)}_{=0} \right] =$

$0 + \alpha [2 + \lambda \max_{a'} \{-1, 0\} - 0] =$
 $0.5 \cdot 2 = 1$



	a_1	a_2
s_1	-1	0
s_2	-1	1
s_3	3	2

Hence, our q-table after five steps is:

$q(s, a):$

	a_1	a_2
s_1	-1	0
s_2	-1	1
s_3	3	2

b, This would correspond to an optimal policy:

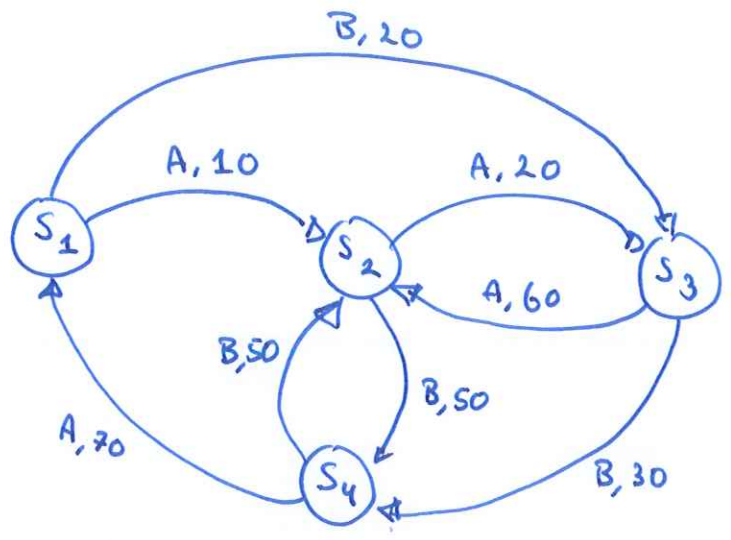
$$a^*(s_1) = \underset{a \in \mathcal{A}}{\text{arg max}} q(s_1, a) = a_2,$$

$$a^*(s_2) = \underset{a \in \mathcal{A}}{\text{arg max}} q(s_2, a) = a_2,$$

$$a^*(s_3) = \underset{a \in \mathcal{A}}{\text{arg max}} q(s_3, a) = a_1.$$

(if we had converged after these iterations)

Ex. 4.7



Consider the deterministic system above with $S = \{s_1, s_2, s_3, s_4\}$, $A = \{A, B\}$ and reward and transitions according to the graph.

Perform Q-learning with $\lambda = 0.9$ and $\alpha = 1$ for the action-sequence $\{A, A, B, A, B, A\}$ initialized in s_1 .

Solution:

Recall the update equation for Q-learning:

$$q(s, a) \leftarrow q(s, a) + \alpha [r(s, a) + \lambda \max_{a'} q(s', a') - q(s, a)]$$

Assume we initialize $q(s, a)$ with zeros:

q:

	A	B
s_1	0	0
s_2	0	0
s_3	0	0
s_4	0	0

• $q(s, a) = q(s, A) \leftarrow$

$$q(s, a) + \alpha [r(s, a) + \lambda \max_{a'} q(s', a') - q(s, a)] =$$

$$\underbrace{q(s, A)}_{=0} + \alpha [\underbrace{r(s, A)}_{=10} + \lambda \max_{a'} \underbrace{q(s', a')}_{=0 \text{ for all } a'} - \underbrace{q(s, A)}_{=0}] =$$

$s = s_1$
 $a = A$
 $s' = s_2$
 $r = 10$

$$0 + 1 [10 + 0.9 \cdot 0 - 0] = 10$$

	A	B
s ₁	10	0
s ₂	0	0
s ₃	0	0
s ₄	0	0

• $q(s_2, A) \leftarrow$

$$q(s_2, A) + \alpha [r(s_2, A) + \lambda \max_{z'} q(s_3, z') - q(s_2, A)] =$$

$\underbrace{q(s_2, A)}_{=0} + \alpha \left[\underbrace{r(s_2, A)}_{=20} + \lambda \max_{z'} \underbrace{q(s_3, z')}_{=0 \text{ for all } z'} - \underbrace{q(s_2, A)}_{=0} \right] =$

$s = s_2$
 $z = A$
 $s' = s_3$
 $r = 20$

$$0 + 1 [20 + 0.9 \cdot 0 - 0] = 20$$

	A	B
s ₁	10	0
s ₂	20	0
s ₃	0	0
s ₄	0	0

• $q(s_3, B) \leftarrow$

$$q(s_3, B) + \alpha [r(s_3, B) + \lambda \max_{z'} q(s_4, z') - q(s_3, B)] =$$

$\underbrace{q(s_3, B)}_{=0} + \alpha \left[\underbrace{r(s_3, B)}_{=30} + \lambda \max_{z'} \underbrace{q(s_4, z')}_{=0 \text{ for all } z'} - \underbrace{q(s_3, B)}_{=0} \right] =$

$s = s_3$
 $z = B$
 $s' = s_4$
 $r = 30$

$$0 + 1 [30 + 0.9 \cdot 0 - 0] = 30$$

	A	B
s ₁	10	0
s ₂	20	0
s ₃	0	30
s ₄	0	0

• $q(s_4, A) \leftarrow$

$$q(s_4, A) + \alpha [r(s_4, A) + \lambda \max_{z'} q(s_1, z') - q(s_4, A)] =$$

$\underbrace{q(s_4, A)}_{=0} + \alpha \left[\underbrace{r(s_4, A)}_{=70} + \lambda \max_{z'} \underbrace{q(s_1, z')}_{=0} - \underbrace{q(s_4, A)}_{=0} \right] =$

$s = s_4$
 $z = A$
 $s' = s_1$
 $r = 70$

$$0 + 1 [70 + \lambda \max_{z'} \{ \overset{A}{10}, \overset{B}{0} \} - 0] =$$

$$70 + 0.9 \cdot 10 = 79$$

	A	B
s ₁	10	0
s ₂	20	0
s ₃	0	30
s ₄	79	0

• $q(s_1, B) \leftarrow$

$$q(s_1, B) + \alpha [r(s_1, B) + \lambda \max_{z'} q(s_3, z') - q(s_1, B)] =$$

$\underbrace{q(s_1, B)}_{=0} + \alpha \left[\underbrace{r(s_1, B)}_{=20} + \lambda \max_{z'} \underbrace{q(s_3, z')}_{=0} - \underbrace{q(s_1, B)}_{=0} \right] =$

$s = s_1$
 $z = B$
 $s' = s_3$
 $r = 20$

$$0 + 1 [20 + \lambda \max_{z'} \{ \overset{A}{0}, \overset{B}{30} \} - 0] =$$

$$20 + 0.9 \cdot 30 = 47$$

	A	B
s_1	10	47
s_2	20	0
s_3	0	30
s_4	79	0

$q(s_3, A) \leftarrow$

$$q(s_3, A) + \alpha [r(s_3, A) + \lambda \max_{a'} q(s_2, a') - q(s_3, A)] =$$

$\underbrace{0}_{=0} \quad \underbrace{60}_{=60} \quad \underbrace{0}_{=0}$

$s = s_3$
 $a = A$
 $s' = s_2$
 $r = 60$

$$0 + 1 [60 + \lambda \max \{20, 0\} - 0] =$$

$$60 + 0.9 \cdot 20 = 78$$

	A	B
s_1	10	47
s_2	20	0
s_3	78	30
s_4	79	0

Hence, our state-action function is

$q(s, a):$

	A	B
s_1	10	47
s_2	20	0
s_3	78	30
s_4	79	0

which corresponds to the greedy policy:

- $a(s_1) = \arg \max_a q(s_1, a) = B,$
- $a(s_2) = A,$
- $a(s_3) = A,$
- $a(s_4) = A.$