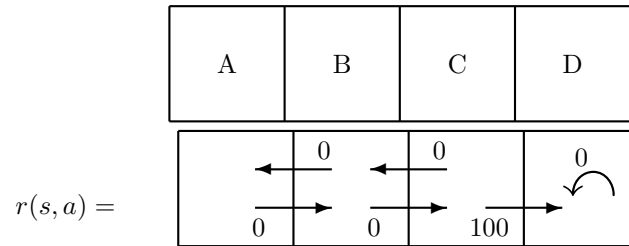


Reinforcement Learning Examples

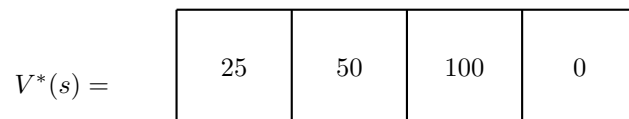
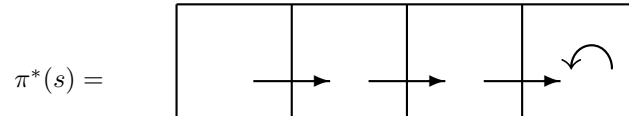
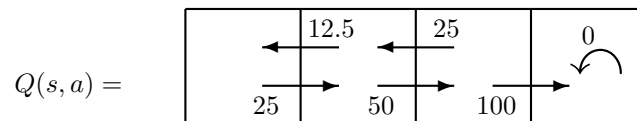
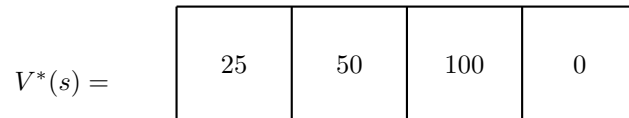
Example 1

A simple board game:



A discount factor: $\gamma = 1/2$.

If we can somehow figure out $V^*(s)$, it is easy to compute Q from V^* , π^* from Q , V^* from π^* .



The offline algorithm: start by guessing \hat{V} .

 $\hat{V} =$

0	0	0	0
---	---	---	---

 $\hat{Q}(s, a) =$

	$\xleftarrow{0}$	$\xleftarrow{0}$	$\xrightarrow{0}$
$\xrightarrow{0}$		$\xrightarrow{0}$	$\xrightarrow{100}$

 $\hat{\pi} =$

$\xleftarrow{\quad}$			$\xrightarrow{\quad}$
$\xrightarrow{\quad}$			$\xrightarrow{\quad}$

 $\hat{V} =$

0	0	100	0
---	---	-----	---

 $\hat{Q}(s, a) =$

	$\xleftarrow{0}$	$\xleftarrow{0}$	$\xrightarrow{0}$
$\xrightarrow{0}$		$\xrightarrow{50}$	$\xrightarrow{100}$

 $\hat{\pi} =$

$\xrightarrow{\quad}$	$\xrightarrow{\quad}$	$\xrightarrow{\quad}$	$\xrightarrow{\quad}$
-----------------------	-----------------------	-----------------------	-----------------------

 $\hat{V} =$

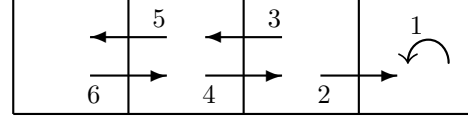
25	50	100	0
----	----	-----	---

The online algorithm.

start with arbitrary Q

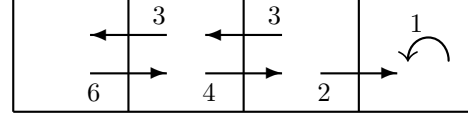
A	B	C	D
---	---	---	---

$$\hat{Q}(s, a) =$$



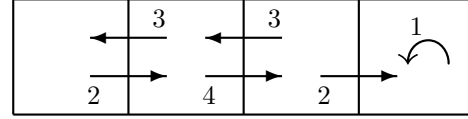
$$Q(B, \text{left}) = 0 + 0.5 \cdot 6 = 3$$

$$\hat{Q}(s, a) =$$



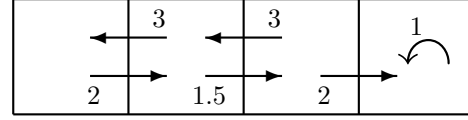
$$Q(A, \text{right}) = 0 + 0.5 \cdot 4 = 2$$

$$\hat{Q}(s, a) =$$



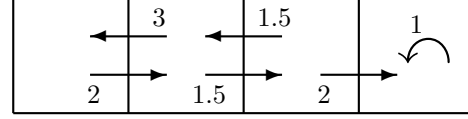
$$Q(B, \text{right}) = 1.5$$

$$\hat{Q}(s, a) =$$



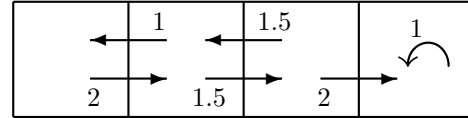
$$Q(C, \text{left}) = 1.5$$

$$\hat{Q}(s, a) =$$



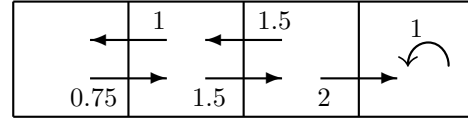
$$Q(B, \text{left}) = 1$$

$$\hat{Q}(s, a) =$$



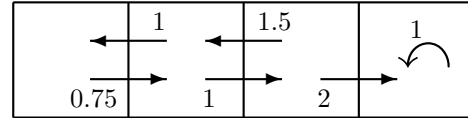
$$Q(A, \text{right}) = 0.75$$

$$\hat{Q}(s, a) =$$



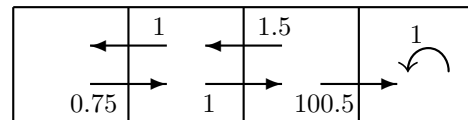
$$Q(B, \text{right}) = 1$$

$$\hat{Q}(s, a) =$$



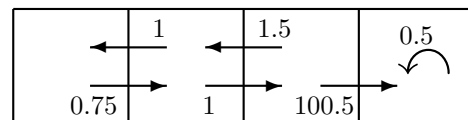
$$Q(C, \text{right}) = 100.5$$

$$\hat{Q}(s, a) =$$



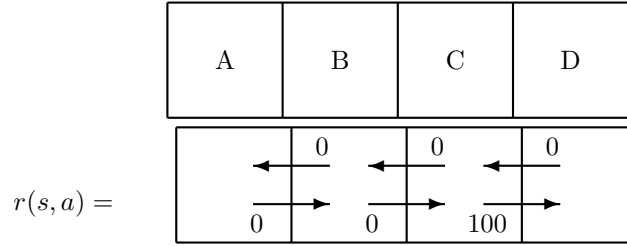
$$Q(D, \text{itself}) = 0.5$$

$$\hat{Q}(s, a) =$$



additional iterations are needed

Example 2



A discount factor: $\gamma = 1/2$.

Use intuition to compute V^* . If we are at D we can go to C , and then back to D and get rewarded.

$$\begin{aligned}
 V^*(D) &= 0 + \frac{1}{2} \cdot 100 + 0 + \frac{1}{2^3} \cdot 100 + \dots \\
 &= 100 \cdot \frac{1}{2} \cdot (1 + (\frac{1}{4}) + (\frac{1}{4})^2 + \dots) \\
 &= 100 \cdot \frac{1}{2} \cdot \frac{1}{1 - 1/4} = 100 \cdot \frac{2}{3} \\
 V^*(C) &= 100 + \frac{1}{2} V^*(D) = 100 \cdot \frac{4}{3} \\
 V^*(B) &= 0 + \frac{1}{2} V^*(C) = 100 \cdot \frac{2}{3} \\
 V^*(A) &= 0 + \frac{1}{2} V^*(B) = 100 \cdot \frac{1}{3}
 \end{aligned}$$

