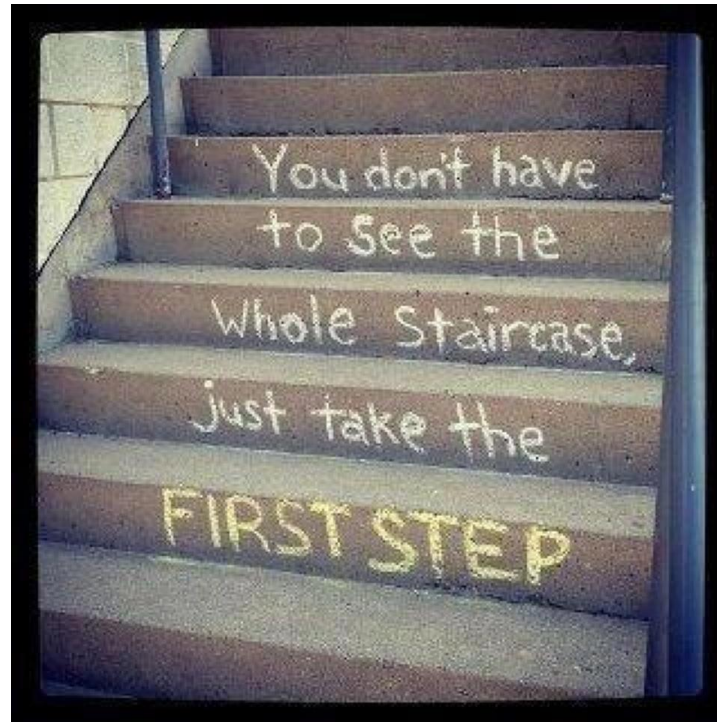


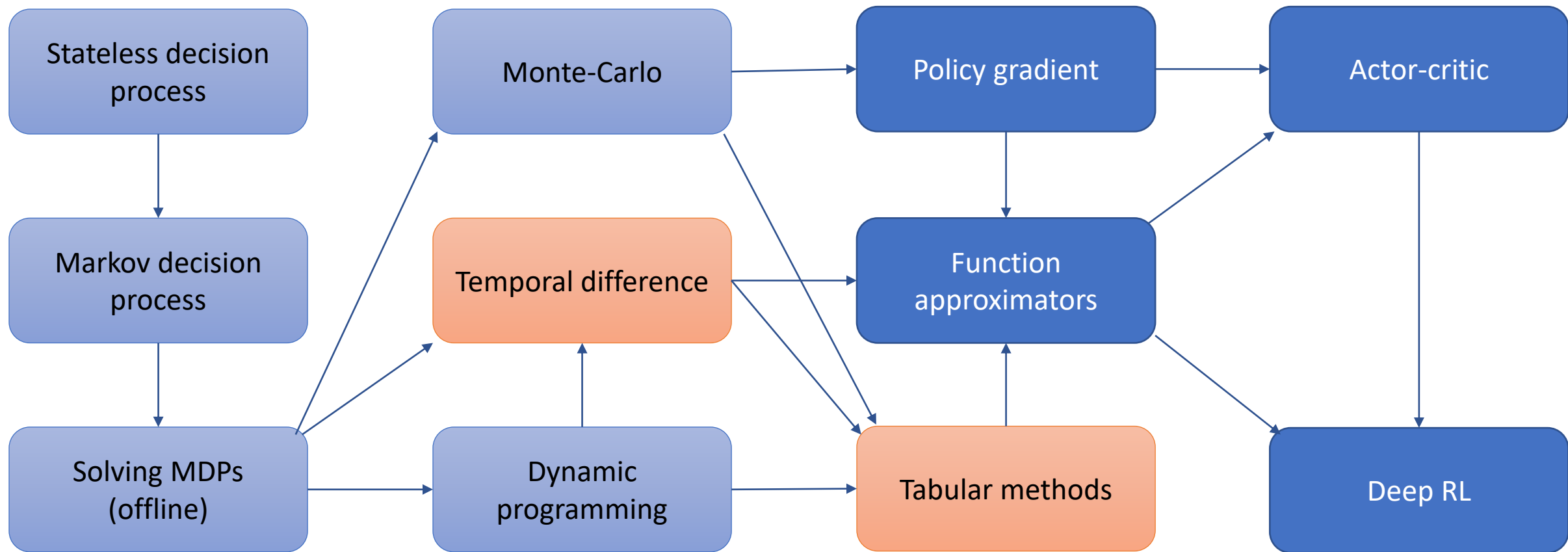
# CSCE-642 Reinforcement Learning

## Chapter 6: Temporal-Difference Learning



Instructor: Guni Sharon

# CSCE-689, Reinforcement Learning



# Solving MDPs so far

## Dynamic programming

- ✓ Off policy
- ✓ local learning, propagating values from neighbors (Bootstrapping)
- ✗ Model based

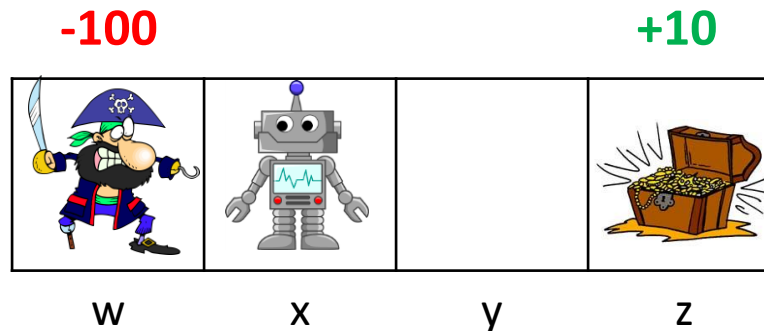
## Monte-Carlo

- ✗ On-policy (though important sampling can be used)
- ✗ Requires a full episode to train on
- ✓ Model free, online learning

- $Q(z, \text{exit}) = 10$
- $Q(y, \rightarrow) = 0 + \gamma \max_a Q(z, a)$
- $Q(x, \rightarrow) = 0 + \gamma \max_a Q(y, a)$

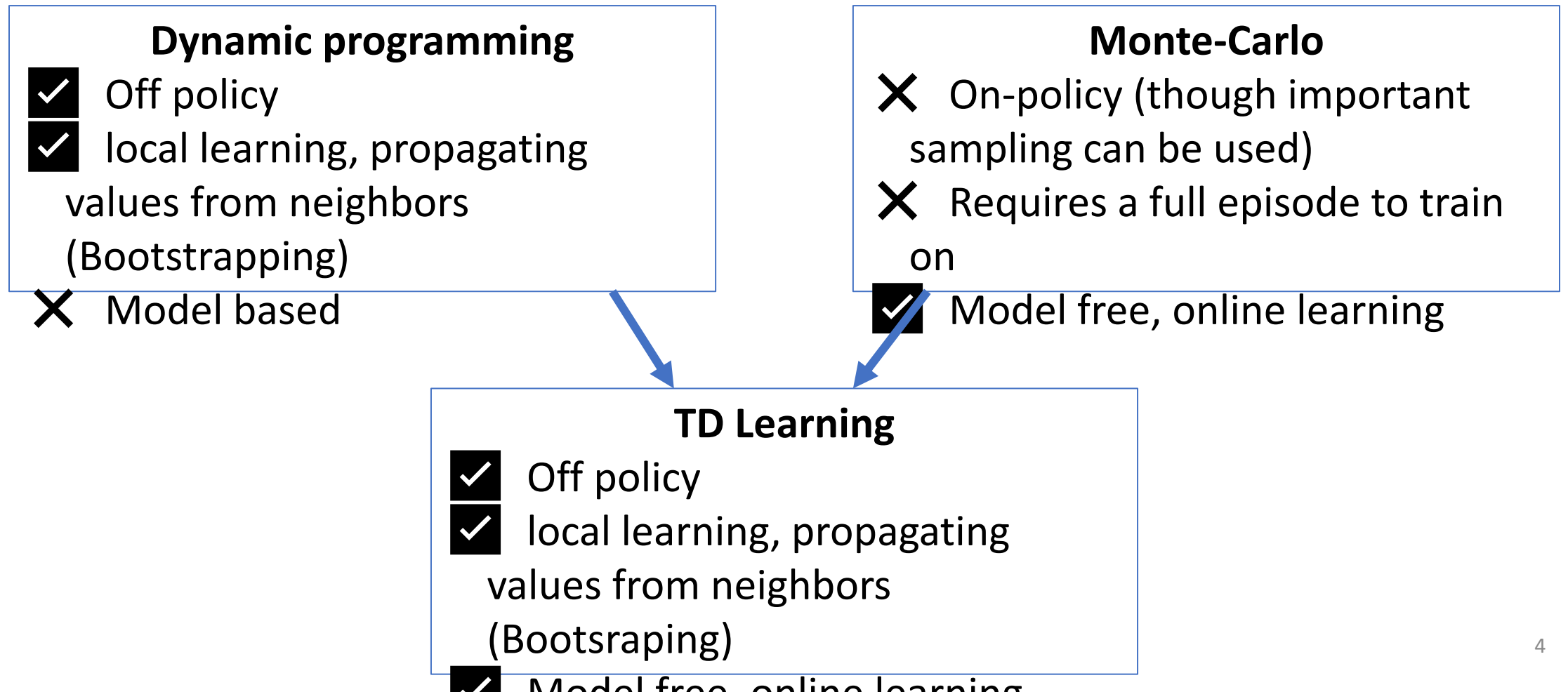
$$q^*(s, a) = \sum_{s'} p(s'|s, a) (r(s, a, s') + \gamma \max_a [q^*(s', a)])$$

$$\gamma = 0.9$$



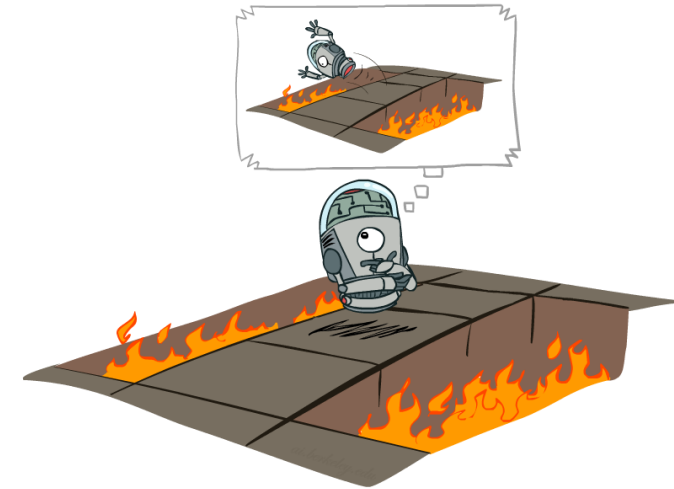
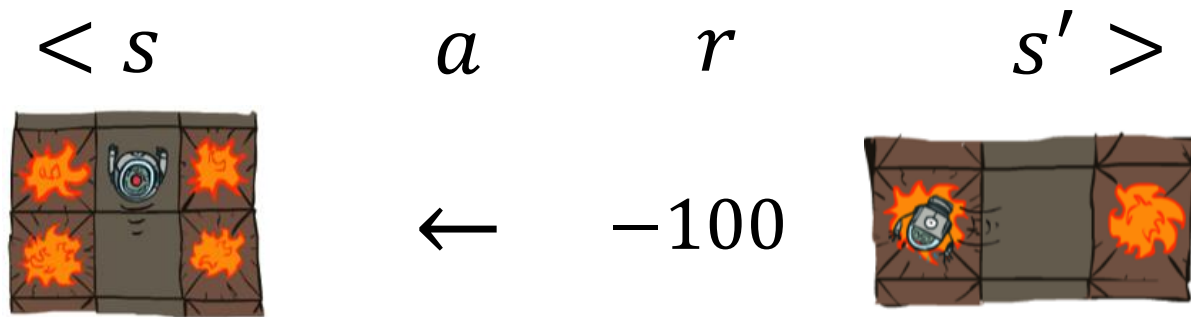
- Episode =  $\{x, y, z, \text{exit}\}$
- $Q(z, \text{exit}) = 10$
- $Q(y, \rightarrow) = 9$
- $Q(x, \rightarrow) = 8.1$

# Fuse DP and MC



# Online Bellman update

- $q^*(s, a) = \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma \max_{a'} [q^*(s', a')])$
- Model and reward function are unknown
- Instead, we observe transitions:  $\langle s_t, a_t, r_{t+1}, s_{t+1} \rangle$



Offline Solution



Online Learning

# Online Bellman update

- $q^*(s, a) = \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma \max_{a'} [q^*(s', a')])$
- Model and reward function are unknown
- Instead, we observe transitions:  $\langle s_t, a_t, r_{t+1}, s_{t+1} \rangle$ 
  - $\mathbb{E}[X] = \sum_{k=1}^{|X|} x_i \Pr\{x_i\} = \frac{\sum_{k=1}^N x_k}{N}$
- $q^*(s, a) = \mathbb{E}_{s' \sim P(s'|s, a)} [\mathbb{E}[R(s, a, s')] + \gamma \max_{a'} [q^*(s', a')]]$
- For a single sample:  $= r_{t+1} + \gamma \max_{a'} [q^*(s_{t+1}, a')]$
- = unbiased estimation of the Bellman update

# Online Bellman update

- $q^*(s, a) = r_{t+1} + \gamma \max_{a'} [q^*(s_{t+1}, a')]$
- Assume that we visited state  $s$  and performed action  $a$  5 times
  - $\langle s, a, 3, x \rangle, \langle s, a, 5, x \rangle, \langle s, a, 1, y \rangle, \langle s, a, 2, y \rangle, \langle s, a, 3, y \rangle$
  - Assume  $\gamma = 0.9$  and  $\max_a [q^*(x, a)] = \max_a [q^*(y, a)] = 10$
- Compute an approximation of  $q^*(s, a)$ :
  - Average of unbiased estimations
  - $((3 + 9) + (5 + 9) + (1 + 9) + (2 + 9) + (3 + 9)) \frac{1}{5}$
- Can  $Q(s, a)$  be updated online?
- Yes, moving average:
  - $Q(s, a) = Q(s, a) + \alpha \left( r_{t+1} + \gamma \max_{a'} [q^*(s', a')] - Q(s, a) \right)$

# Temporal difference learning

- $Q(s, a) = Q(s, a) + \alpha \left( R_{t+1} + \gamma \max_{a'} [q^*(s', a')] - Q(s, a) \right)$
- But we don't know  $q^*(s', a')$
- Use the estimation  $Q(s', a)$
- $Q(s, a) = Q(s, a) + \alpha \left( \underbrace{R_{t+1} + \gamma \max_{a'} [Q(s', a')] - Q(s, a)}_{\text{Temporal difference error: } \delta} \right)$



# Temporal difference learning

- $Q(s, a) = Q(s, a) + \alpha \left( R_{t+1} + \gamma \max_{a'} [Q(s', a')] - Q(s, a) \right)$ 
  - $V(s) = V(s) + \alpha (R_{t+1} + \gamma V(s') - V(s))$
- Update estimate based on other estimates
- Model free
- Online, incremental learning
- Guaranteed to converge to the true value!
  - Some conditions on the step size,  $\alpha$  (see slide #19 in 2Multi\_armed\_bandits.pptx)
- Usually converges faster than MC methods

# SARSA: On-policy TD Control

- $Q(s, a) = Q(s, a) + \alpha \left( R_{t+1} + \gamma \max_{a'} [Q(s', a')] - Q(s, a) \right)$
- Replace  $\max_{a'} [Q(s', a')]$  with values from the observed transition
  - $\langle s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1} \rangle$
- $Q(s_t, a_t) = Q(s_t, a_t) + \alpha (r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$
- SARSA converges with probability 1 to an optimal policy and action-values as long as all state–action pairs are visited infinitely often, and the policy converges in the limit to the greedy policy
  - Which can be arranged, for example, with  $\epsilon$ -greedy policies by setting  $\epsilon = 1/t$

# SARSA: On-policy TD Control

## Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

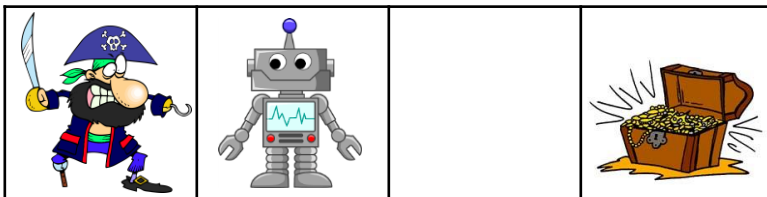
$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal

-100

$\gamma = 0.9$

+10



w

x

y

z

# SARSA: On-policy TD Control

## Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

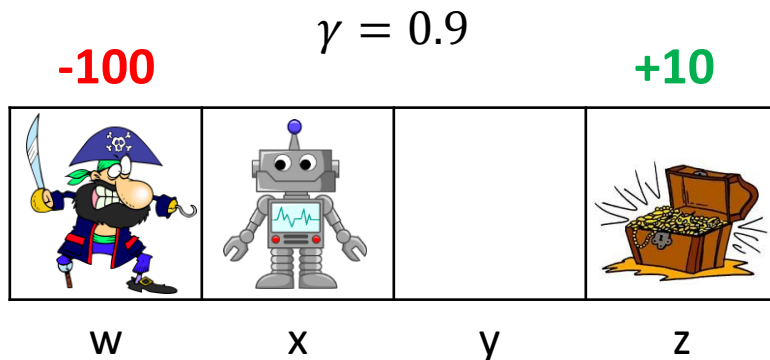
Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal



$Q =$

0	0,0	0,0	0
w	x	y	z

# SARSA: On-policy TD Control

## Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

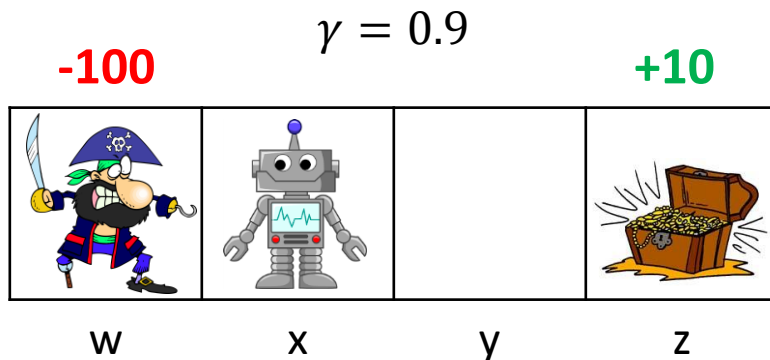
Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal



$x$	$\leftarrow$	$0$	$w$	$null$
$S$	$A$	$R$	$S'$	$A'$

$Q =$

$0$	$0,0$	$0,0$	$0$
$w$	$x$	$y$	$z$

# SARSA: On-policy TD Control

## Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

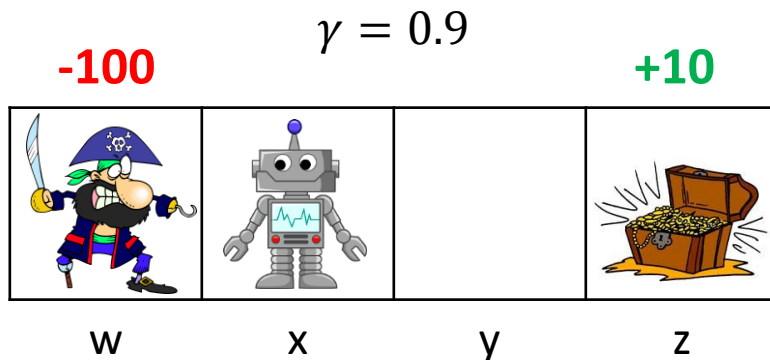
Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal



$x$	$\leftarrow$	$0$	$w$	$exit$
$S$	$A$	$R$	$S'$	$A'$

$Q =$

$0$	$0,0$	$0,0$	$0$
$w$	$x$	$y$	$z$

# SARSA: On-policy TD Control

## Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

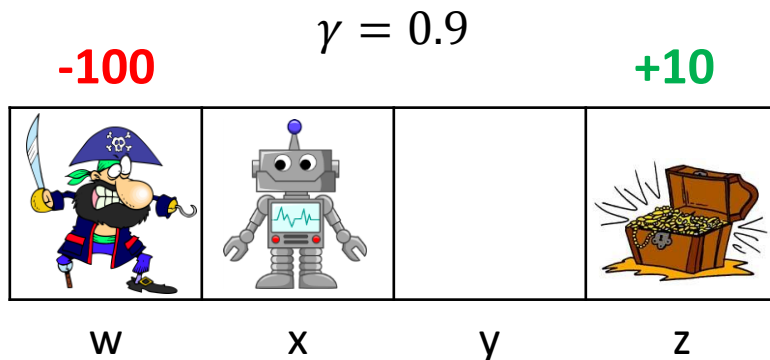
Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal



$x$	$\leftarrow$	$0$	$w$	$exit$
$S$	$A$	$R$	$S'$	$A'$

$Q =$

$0$	<b><math>0,0</math></b>	$0,0$	$0$
$w$	$x$	$y$	$z$

# SARSA: On-policy TD Control

## Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

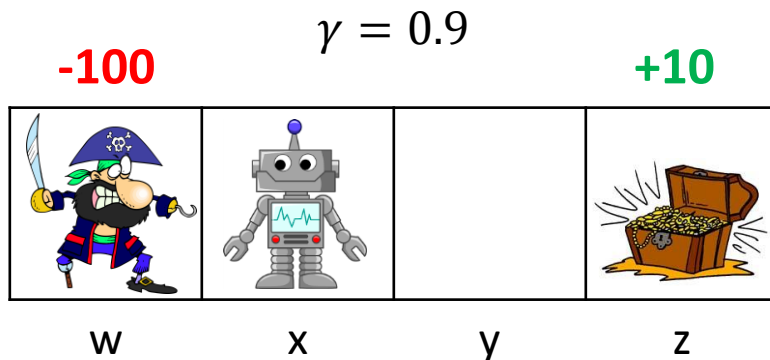
Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal



w	exit	0	w	exit
S	A	R	S'	A'

$Q =$

0	0,0	0,0	0
w	x	y	z



# SARSA: On-policy TD Control

## Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

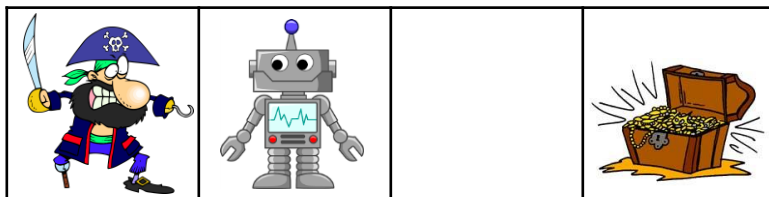
$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal

$\gamma = 0.9$

-100

+10



w

x

y

z

w	exit	-100	ter	exit
S	A	R	S'	A'

$Q =$

0	0,0	0,0	0
w	x	y	z

# SARSA: On-policy TD Control

## Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

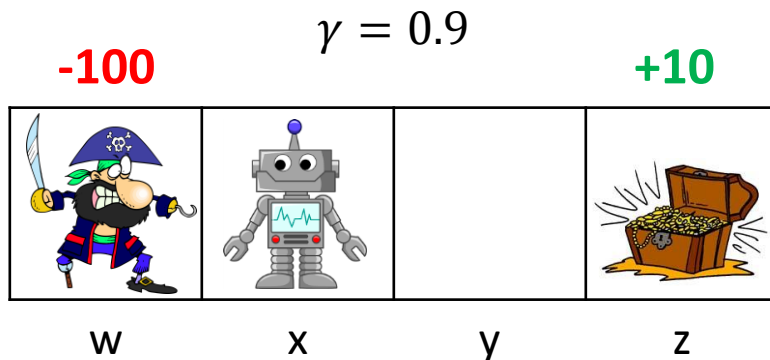
Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal



w	exit	-100	ter	.
S	A	R	S'	A'

$Q =$

-100	0,0	0,0	0
w	x	y	z

# SARSA: On-policy TD Control

## Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

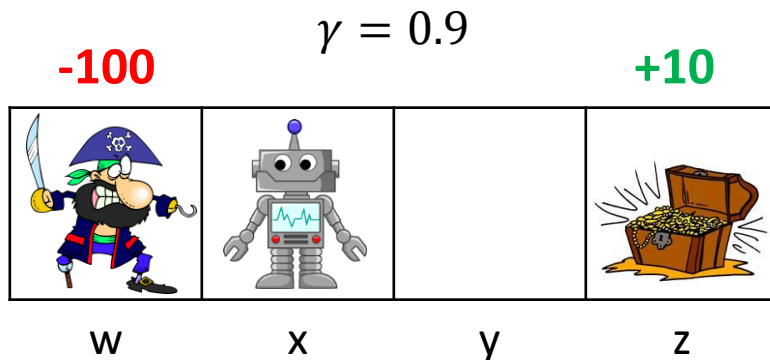
Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal



$x$	$\leftarrow$	$-100$	$ter$	$\cdot$
$S$	$A$	$R$	$S'$	$A'$

$Q =$

$-100$	$0,0$	$0,0$	$0$
w	x	y	z

# SARSA: On-policy TD Control

## Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

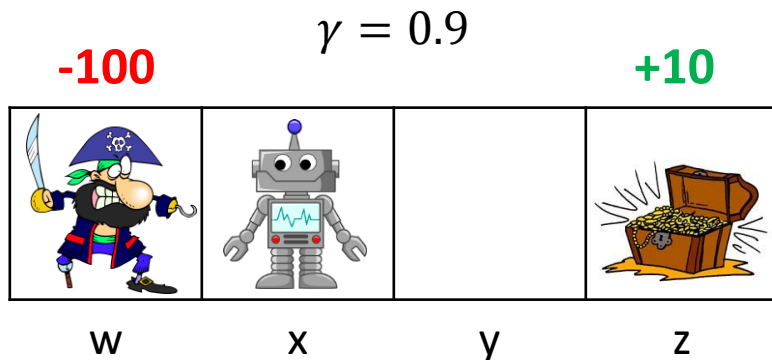
Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal



$x$	$\leftarrow$	$0$	$w$	$exit$
$S$	$A$	$R$	$S'$	$A'$

$Q =$

<b>-100</b>	<b>0,0</b>	<b>0,0</b>	<b>0</b>
w	x	y	z

# SARSA: On-policy TD Control

## Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

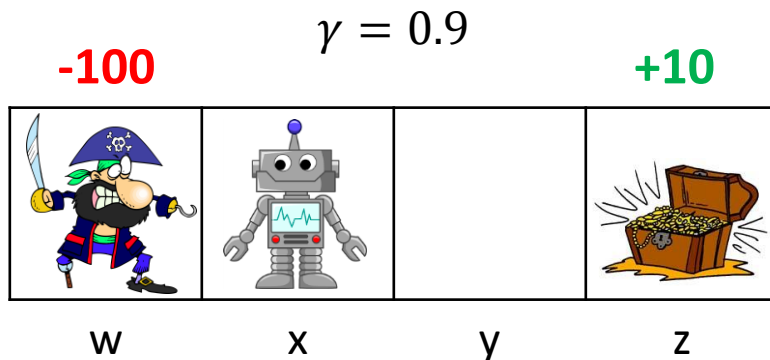
Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal



$x$	$\leftarrow$	$0$	$w$	$exit$
$S$	$A$	$R$	$S'$	$A'$

$Q =$

-100	-90,0	0,0	0
w	x	y	z

# SARSA: On-policy TD Control

## Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A';$

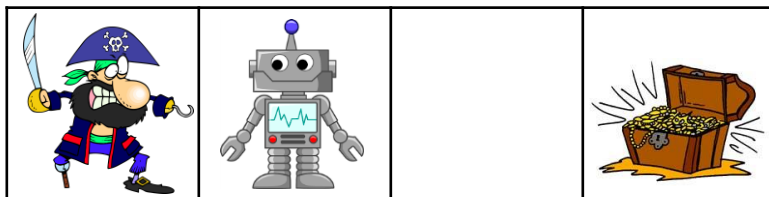
until  $S$  is terminal

And so on...

$\gamma = 0.9$

-100

+10



w

x

y

z

w	exit	0	w	exit
S	A	R	S'	A'

$Q =$

-100	-90,0	0,0	0
w	x	y	z

# Q-learning: Off-policy TD Control

- Use the original TD update rule
- $Q(s, a) = Q(s, a) + \alpha \left( R_{t+1} + \gamma \max_{a'} [Q(s', a')] - Q(s, a) \right)$
- Converge on the state-action value for the optimal policy, i.e.,  $q^*$ 
  - Assuming that every state-action pair is visited infinitely often
- Follows from the proof of convergence for the Bellman function
  - See slides #25,26 in “3MDPs+DP.pptx”

# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

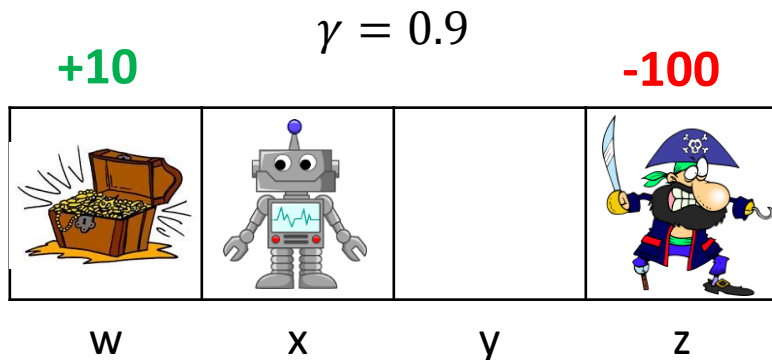
Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$

until  $S$  is terminal



$x$	<i>null</i>	<i>null</i>	<i>null</i>
$S$	$A$	$R$	$S'$

$Q =$

0	0,0	0,0	0
w	x	y	z



# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

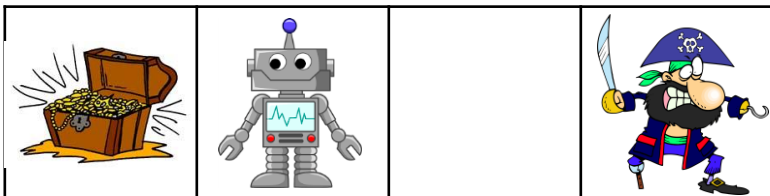
$S \leftarrow S'$

until  $S$  is terminal

+10

$\gamma = 0.9$

-100



w

x

y

z

$x$	$\rightarrow$	$0$	$y$
$S$	$A$	$R$	$S'$

$Q =$

$0$	$0,0$	$0,0$	$0$
$w$	$x$	$y$	$z$

# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

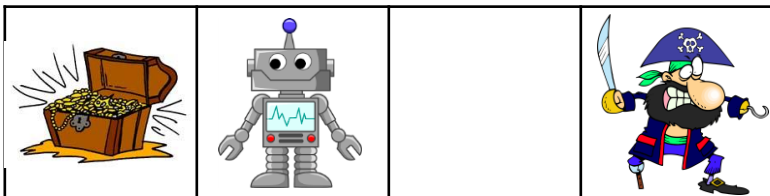
$S \leftarrow S'$

until  $S$  is terminal

+10

$\gamma = 0.9$

-100



w

x

y

z

$x$	$\rightarrow$	$0$	$y$
$S$	$A$	$R$	$S'$

$Q =$

$0$	$0, 0$	$0, 0$	$0$
w	x	y	z

# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

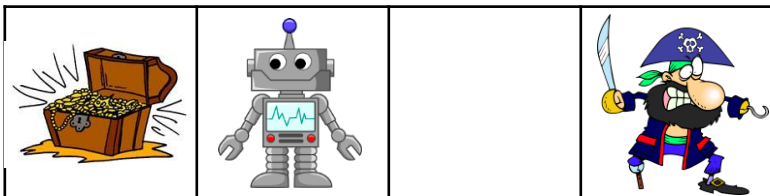
$S \leftarrow S'$

until  $S$  is terminal

+10

$\gamma = 0.9$

-100



w

x

y

z

$y$	$\rightarrow$	0	$y$
$S$	$A$	$R$	$S'$

$Q =$

0	0,0	0,0	0
w	x	y	z

# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

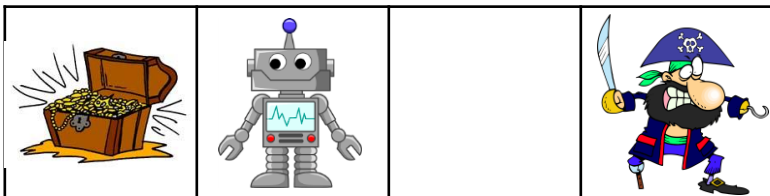
$S \leftarrow S'$

until  $S$  is terminal

+10

$\gamma = 0.9$

-100



w

x

y

z

y	→	0	z
S	A	R	S'

Q =	0	0,0	0,0	0
	w	x	y	z

# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

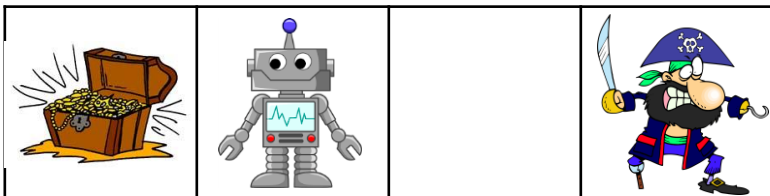
$S \leftarrow S'$

until  $S$  is terminal

+10

$\gamma = 0.9$

-100



w

x

y

z

z	→	0	z
S	A	R	S'

Q =	0	0,0	0,0	0
	w	x	y	z

# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

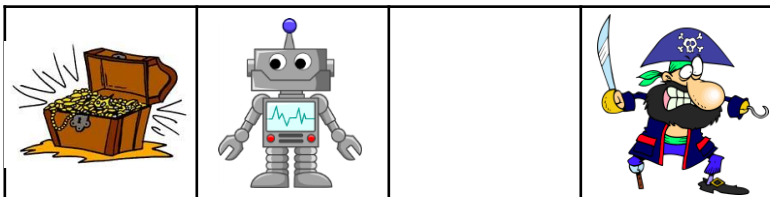
$S \leftarrow S'$

until  $S$  is terminal

+10

$\gamma = 0.9$

-100



w

x

y

z

z	exit	-100	.
---	------	------	---

S

A

R

S'

$Q =$

0	0,0	0,0	0
---	-----	-----	---

w

x

y

z

# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

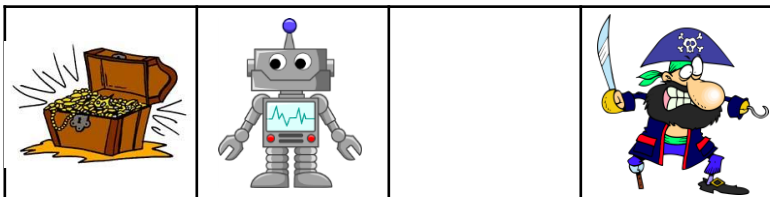
$S \leftarrow S'$

until  $S$  is terminal

+10

$\gamma = 0.9$

-100



w

x

y

z

z	exit	-100	.
---	------	------	---

S

A

R

S'

$Q =$

0	0,0	0,0	-100
---	-----	-----	------

w

x

y

z

# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

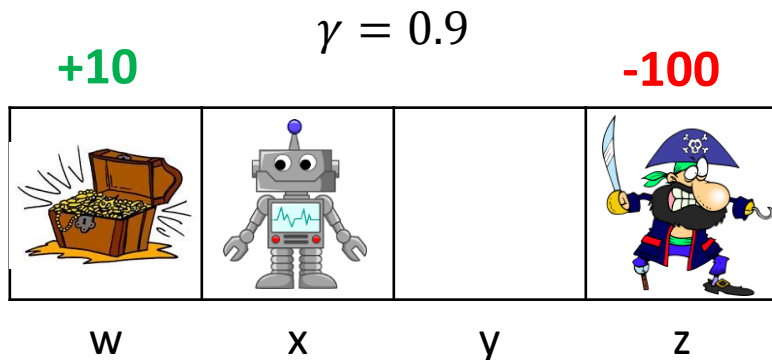
Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$

until  $S$  is terminal



$x$	$exit$	$-100$	$\cdot$
$S$	$A$	$R$	$S'$

$Q =$

0	0,0	0,0	-100
w	x	y	z



# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

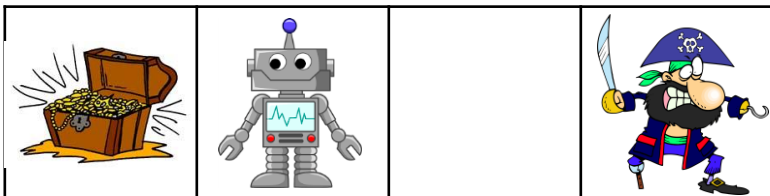
$S \leftarrow S'$

until  $S$  is terminal

+10

$\gamma = 0.9$

-100



w

x

y

z

$x$	$\rightarrow$	$0$	$y$
$S$	$A$	$R$	$S'$

$Q =$

$0$	$0,0$	$0,0$	$-100$
$w$	$x$	$y$	$z$

# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

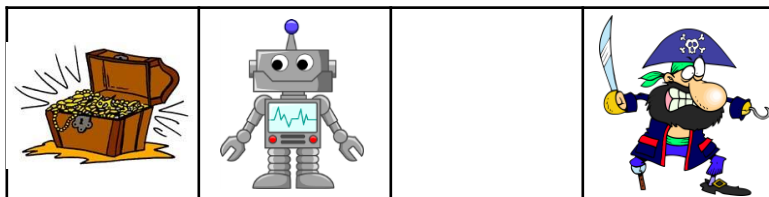
$S \leftarrow S'$

until  $S$  is terminal

+10

$\gamma = 0.9$

-100



w

x

y

z

$y$	$\rightarrow$	0	$z$
$S$	$A$	$R$	$S'$

$Q =$

0	0,0	0,0	-100
w	x	y	z

# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

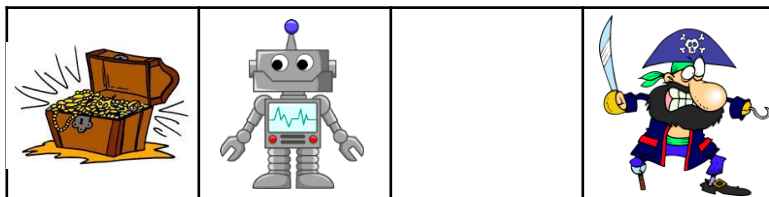
$S \leftarrow S'$

until  $S$  is terminal

+10

$\gamma = 0.9$

-100



w

x

y

z

$y$	$\rightarrow$	0	$z$
$S$	$A$	$R$	$S'$

$Q =$

0	0,0	0,-90	-100
w	x	y	z

# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

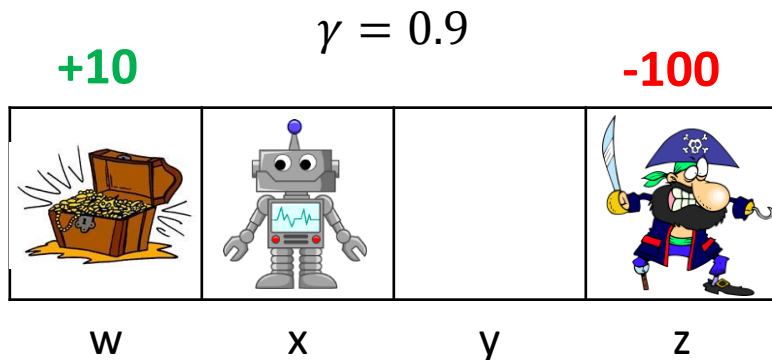
Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$

until  $S$  is terminal



$x$	$\rightarrow$	$0$	$z$
$S$	$A$	$R$	$S'$

$Q =$

$0$	$0,0$	$0,-90$	$-100$
$w$	$x$	$y$	$z$

# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

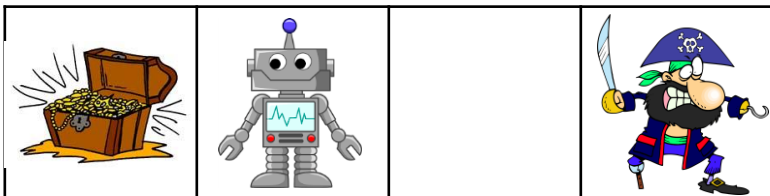
$S \leftarrow S'$

until  $S$  is terminal

+10

$\gamma = 0.9$

-100



w

x

y

z

$x$	$\rightarrow$	$0$	$y$
$S$	$A$	$R$	$S'$

$Q =$

$0$	$0,0$	$0,-90$	$-100$
$w$	$x$	$y$	$z$

# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

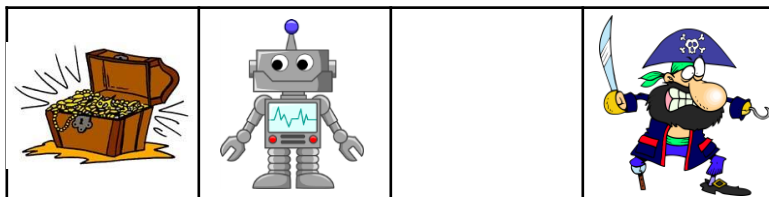
$S \leftarrow S'$

until  $S$  is terminal

+10

$\gamma = 0.9$

-100



w

x

y

z

$x$	$\rightarrow$	$0$	$y$
$S$	$A$	$R$	$S'$

$Q =$

$0$	$0, 0$	$0, -90$	$-100$
$w$	$x$	$y$	$z$

# Q-learning: Off-policy TD Control

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$

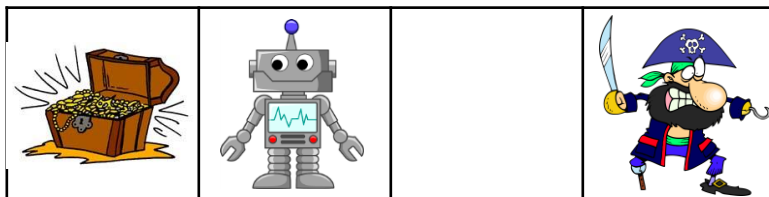
until  $S$  is terminal

And so on...

+10

$\gamma = 0.9$

-100



w

x

y

z

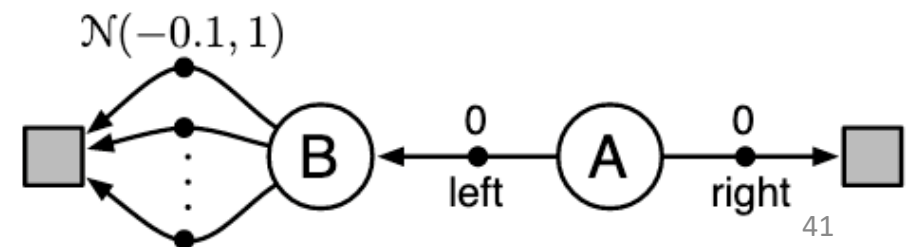
$x$	$\rightarrow$	$0$	$y$
$S$	$A$	$R$	$S'$

$Q =$

$0$	$0, 0$	$0, -90$	$-100$
w	x	y	z

# Maximization bias in Q-learning

- Consider an MDP with two non-terminal states A and B
- Episodes start in A with a choice between two actions, left and right
- The right action transitions immediately to the terminal state with a reward and return of zero
- The left action transitions to B, also with a reward of zero
- From B there are many possible actions all of which cause immediate termination with a reward drawn from a normal distribution with mean  $-0.1$  and variance  $1.0$
- Thus, the expected return for any trajectory starting with left is  $-0.1$ , and thus taking left in state A is a mistake





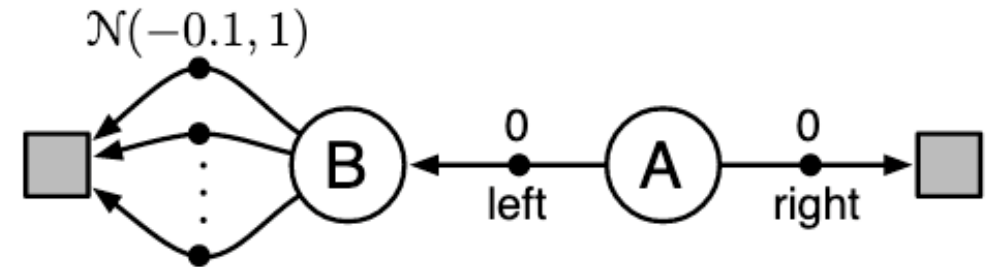
# Maximization bias in Q-learning

- Assume we observed 4 episodes:

- $\{A, \rightarrow, 0\}$
- $\{A, \leftarrow, 0, B, a_1, -1.1\}$
- $\{A, \leftarrow, 0, B, a_2, 0.9\}$
- $\{A, \leftarrow, 0, B, a_3, -0.1\}$

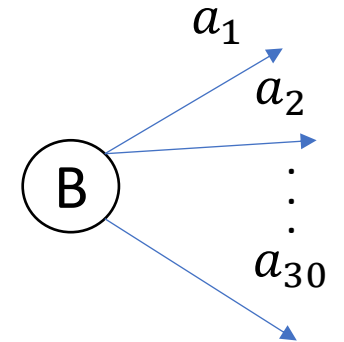
- $\max_a [Q(B, a)] = 0.9$  yet  $q^*(B, \cdot) = -0.1$

- Maximization bias is common when action outcomes are noisy
- Some actions will be evaluated following a lower-than-expected sampled reward and some with a higher-than-expected sampled reward
- The max operator biases towards higher-than-expected sampled reward



# Double Q-learning

- Assume visiting state B a total of 60 times (2 times per action)
- $a_i \sim \mathcal{N}(-0.1, 1)$
- $\mathbb{E} \left( \max_a [Q(B, a)] \right) = \sim^* 1.3 \neq -0.1$
- How can we fix this bias?
  - Store and update 2 independent Q tables:  $Q_1, Q_2$
  - For each observed transition update one Q table (and not the other!)
  - Use one for choosing maximizing action and the other for retrieving the value
- $\mathbb{E} \left( Q_1(B, \arg\max_a [Q_2(B, a)]) \right) = -0.1$



# Double Q learning is unbiased

- We are actually interested in the action that maximizes the expected reward
  - $\max_a [\mathbb{E}[Q(s, a)]]$
- Initially  $\mathbb{E} \left( \max_a [Q(s, a)] \right) \neq \max_a [\mathbb{E}[Q(s, a)]]$ 
  - Due to the maximization bias
- Double Q learning  $\mathbb{E} \left( Q_1(s, \operatorname{argmax}_a [Q_2(s, a)]) \right) = \max_a [\mathbb{E}[Q(s, a)]]$ 
  - Because the returned value (from  $Q_1$ ) is independent of the maximizing action (from  $Q_2$ )

# Double Q-learning

## Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$

until  $S$  is terminal

## Double Q-learning

Initialize  $Q_1(s, a)$  and  $Q_2(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily

Initialize  $Q_1(\text{terminal-state}, \cdot) = Q_2(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q_1$  and  $Q_2$  (e.g.,  $\epsilon$ -greedy in  $Q_1 + Q_2$ )

Take action  $A$ , observe  $R, S'$

With 0.5 probability:

$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha (R + \gamma Q_2(S', \arg \max_a Q_1(S', a)) - Q_1(S, A))$

else:

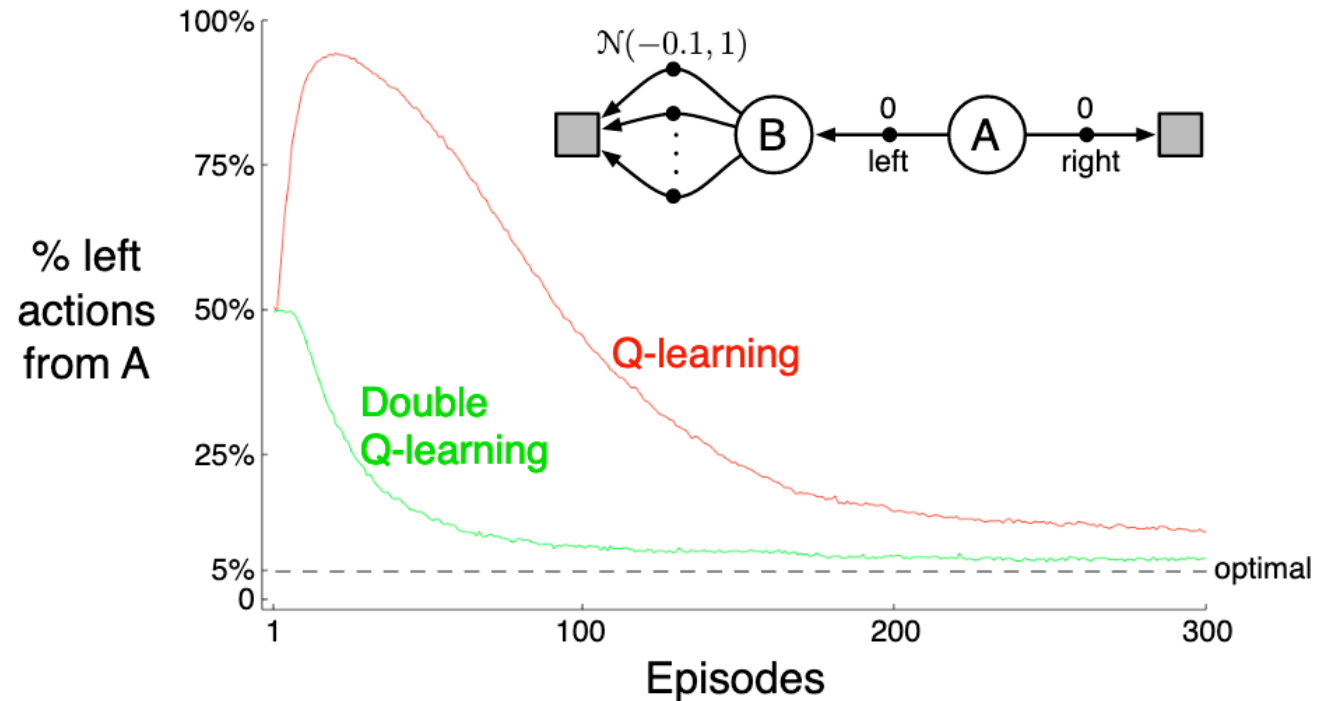
$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha (R + \gamma Q_1(S', \arg \max_a Q_2(S', a)) - Q_2(S, A))$

$S \leftarrow S'$

until  $S$  is terminal

# Double Q-learning

- From B there are many possible actions all of which cause immediate termination with a reward drawn from  $\mathcal{N}(-0.1, 1)$
- Taking left in state A is always a mistake



# What did we learn?

- Temporal difference = online computation of the TD error,  $\delta$
- Allows us to perform unbiased online Bellman updates without any knowledge of the model
- In many cases converges faster than MC (no need to complete an episode)
- SARSA (on policy) provides unbiased TD learning through on policy estimations
- Q learning (off policy) might suffer from a maximization bias that can be addressed with double Q learning
- Both SARSA and Q learning are guaranteed to converge to the optimal state/action values in the tabular case + appropriate learning rate and epsilon decay

# What next?

- **Class:**  $n$ -step Bootstrapping
- **Assignments:**
  - Tabular Q-Learning
  - SARSA
  - Due by October 7, EOD, through Canvas
- **Quiz (on Canvas):**
  - TD learning
  - By Sep 18, EOD
- **Project:**
  - Define your research hypothesis/question