



BITS Pilani
Pilani | Dubai | Goa | Hyderabad

Deep Reinforcement Learning
2023-24 Second Semester, M.Tech (AIML)

Session #6-7-8: Monte Carlo Methods

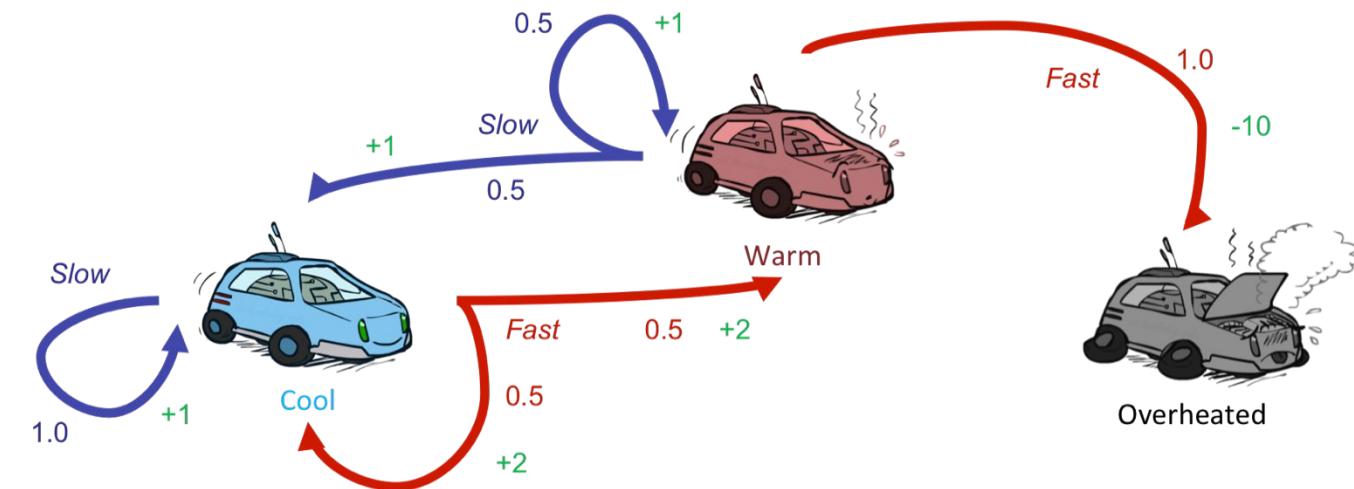


Agenda for the class

- Introduction
- On-Policy Monte Carlo Methods
- Off-Policy Monte Carlo Methods

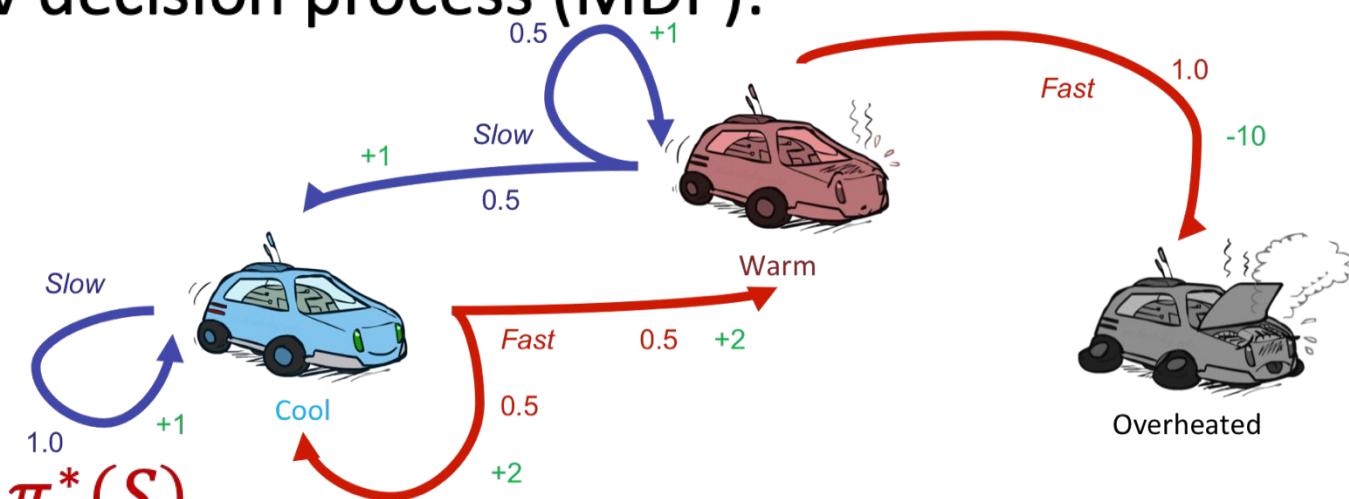
Introduction

- Recollect the problem
 - We need to learn a policy that takes us as far and as faster possible;



Introduction

- Still assume an underlying Markov decision process (MDP):
 - A set of states $s \in S$
 - A set of actions A
 - A model $P(s'|s, a)$
 - A reward function $R(s, a, s')$
 - A discount factor γ
 - Still looking for the best policy $\pi^*(S)$



Introduction

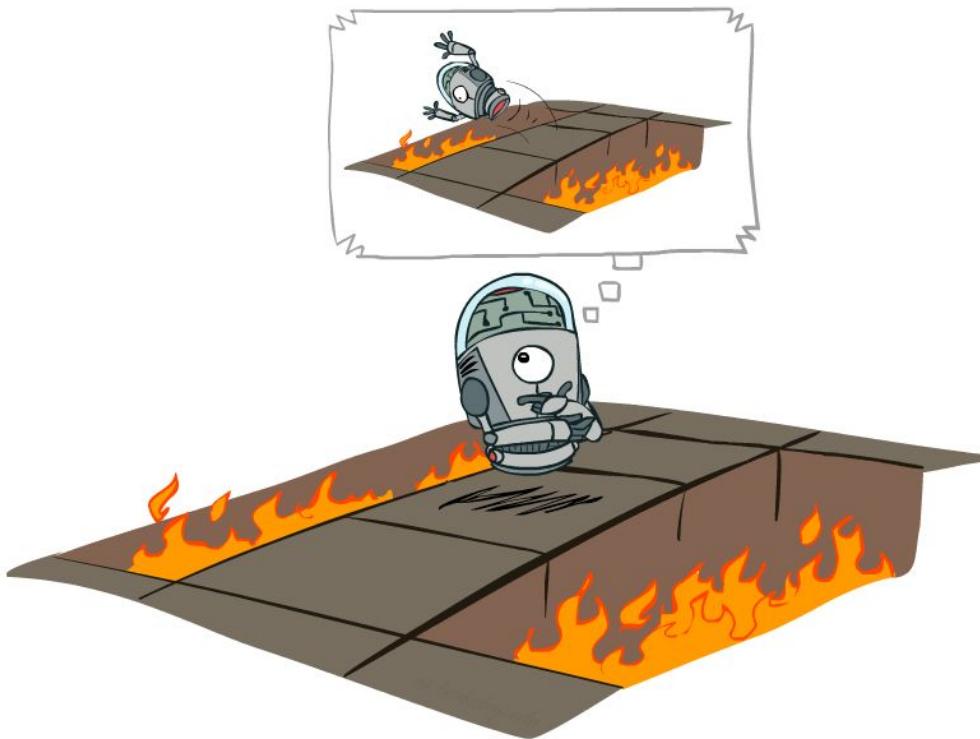
- Still assume an underlying Markov decision process (MDP):

- A **set of states** $s \in S$
- A **set of actions** A
- A **model** $P(s'|s, a)$
- A **reward function** $R(s, a, s')$
- A discount factor γ
- Still looking for the best policy $\pi^*(S)$



- New twist: **don't know the model and the reward function**
 - That is, we don't know the actions' outcome
 - Must interact with the environment to learn

(Aside) Offline vs. Online (RL)



Offline Optimization



Online Learning



Monte Carlo Methods

- Monte Carlo methods are a **broad class of computational algorithms** that *rely on repeated random sampling to obtain numerical results*
- The underlying concept is to obtain unbiased samples from a complex/unknown distribution through a random process
- They are often used in physical and mathematical problems and are most useful when it is difficult or impossible to compute a solution analytically
 - Weather prediction
 - Computational biology
 - Computer graphics
 - Finance and business
 - Sport game prediction



First-visit Monte-Carlo Policy Evaluation

[estimate $V\pi(s)$]

Initialize:

$\pi \leftarrow$ policy to be evaluated

$V \leftarrow$ an arbitrary state-value function

$Returns(s) \leftarrow$ an empty list, for all $s \in \mathcal{S}$

Repeat forever:

(a) Generate an episode using π

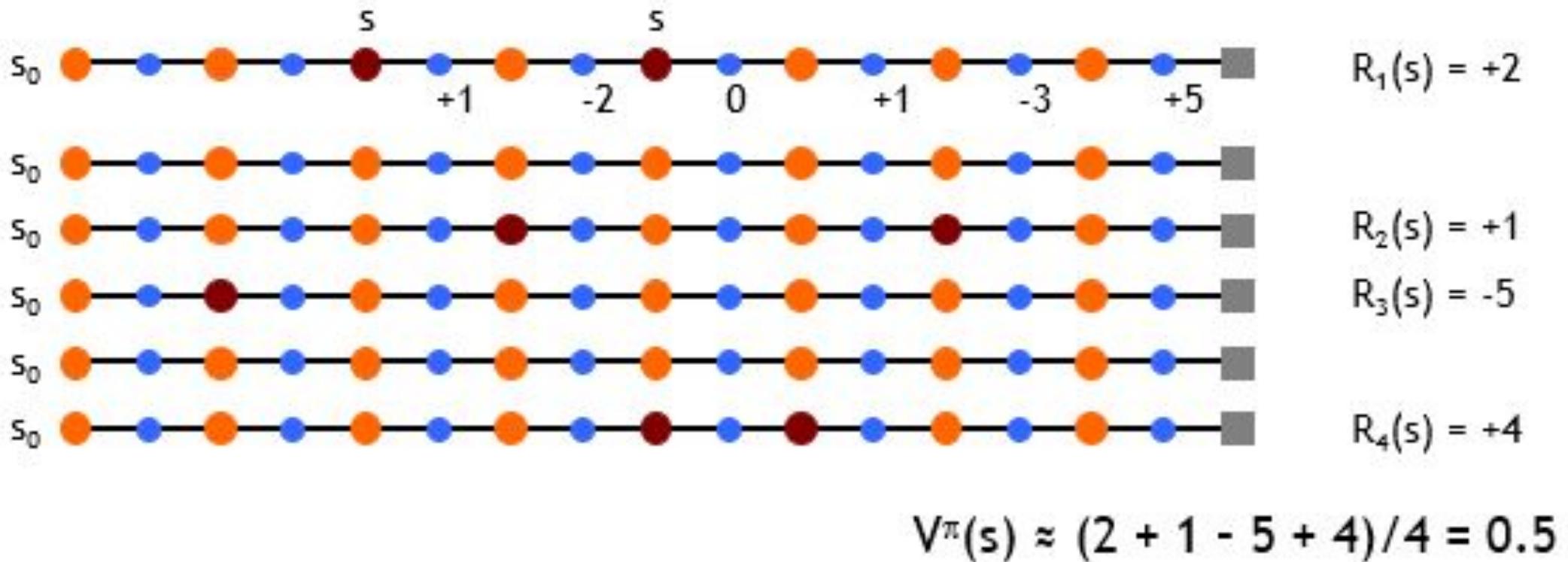
(b) For each state s appearing in the episode:

$R \leftarrow$ return following the first occurrence of s

Append R to $Returns(s)$

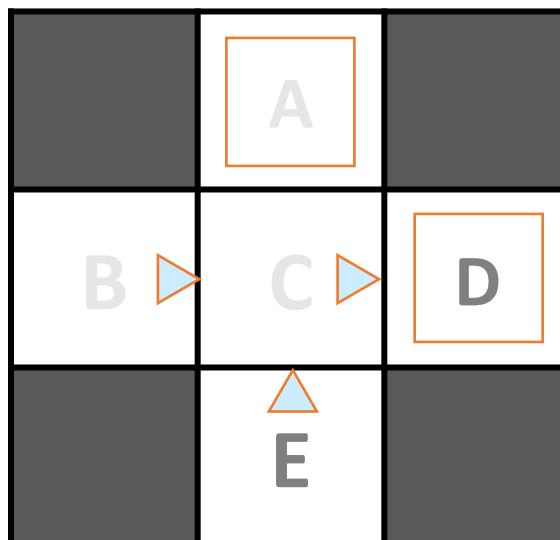
$V(s) \leftarrow \text{average}(Returns(s))$

Ex-1:First-visit Monte-Carlo Policy Evaluation [estimate $V\pi(s)$]



Ex-2: First-visit Monte-Carlo Policy Evaluation [estimate $V\pi(s)$]

Input Policy π



Assume: $\gamma = 1$

Observed Episodes (Training)

Episode 1

B, east, C, -1
C, east, D, -1
D, exit, , +10

Episode 2

B, east, C, -1
C, east, D, -1
D, exit, , +10

Episode 3

E, north, C, -1
C, east, D, -1
D, exit, , +10

Episode 4

E, north, C, -1
C, east, A, -1
A, exit, , -10

Output Values

	-10	
A	+8	+4
B	C	D
	-2	
E		

Problems with MC Evaluation

- What's good about direct evaluation?
 - It's easy to understand
 - It doesn't require any knowledge of the underlying model
 - It converges to the true expected values
- What bad about it?
 - It wastes information about transition probabilities
 - Each state must be learned separately
 - So, it takes a long time to learn

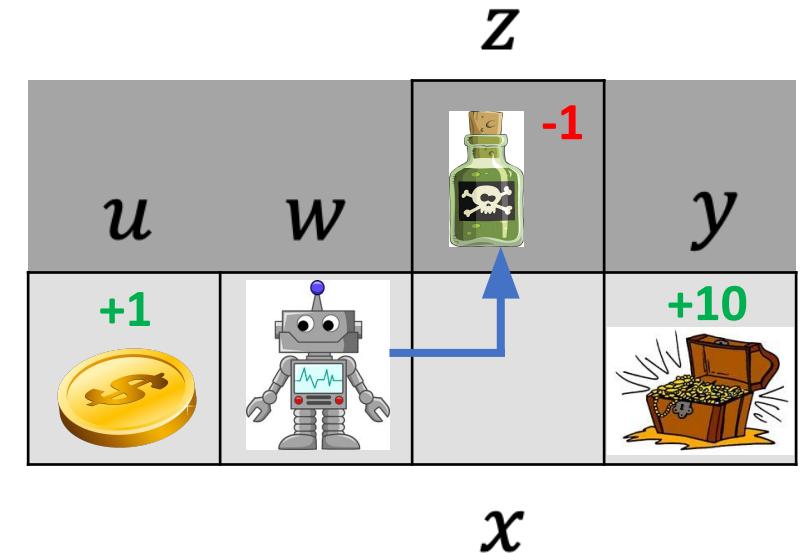
Output Values

	-10 A	
+8 B	+4 C	+10 D
	-2 E	

Think : If B and E both go to C with the same probability, how can their values be different?

What about exploration?

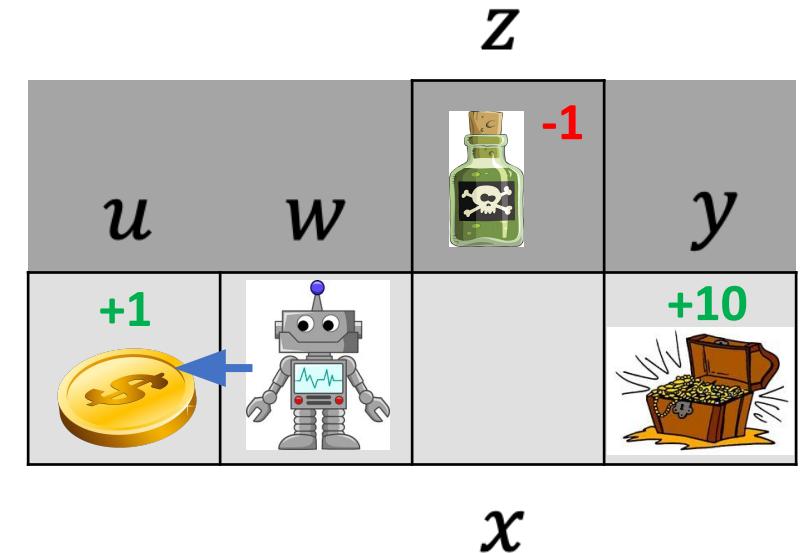
- $S = \{u, w, x, y, z\}$
- $A = \{N, E, S, W, \text{exit}\}$
- Reward:
 - $r(u, \text{exit}) = 1$
 - $r(z, \text{exit}) = -1$
 - $r(y, \text{exit}) = 10$



State	$\pi(s)$	$Q_\pi(N)$	$Q_\pi(E)$	$Q_\pi(S)$	$Q_\pi(W)$	$Q_\pi(\text{exit})$
u	exit	NA	NA	NA	NA	0
w	E	0	0 -1	0	0	NA
x	N	0 -1	0	0	0	NA
y	exit	NA	NA	NA	NA	0
z	exit	NA	NA	NA	NA	0 -1

What about exploration?

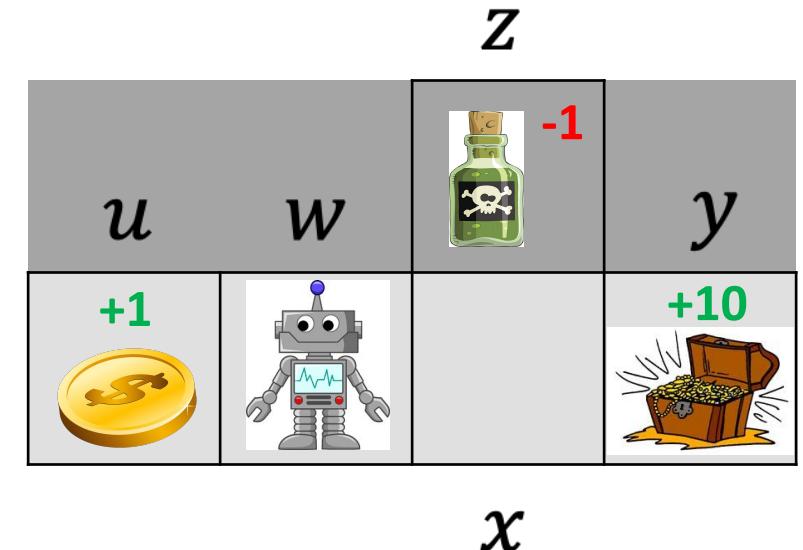
- $S = \{u, w, x, y, z\}$
- $A = \{N, E, S, W, \text{exit}\}$
- Reward:
 - $r(u, \text{exit}) = 1$
 - $r(z, \text{exit}) = -1$
 - $r(y, \text{exit}) = 10$



State	$\pi(s)$	$Q_\pi(N)$	$Q_\pi(E)$	$Q_\pi(S)$	$Q_\pi(W)$	$Q_\pi(\text{exit})$
u	exit	NA	NA	NA	NA	0, 1
w	N, W	0	-1	0	0, 1	NA
x	N, E	-1	0	0	0	NA
y	exit	NA	NA	NA	NA	0
z	exit	NA	NA	NA	NA	-1

What about exploration?

- $S = \{u, w, x, y, z\}$
- $A = \{N, E, S, W, \text{exit}\}$
- Reward:
 - $r(u, \text{exit}) = 1$
 - $r(z, \text{exit}) = -1$
 - $r(y, \text{exit}) = 10$



State	$\pi(s)$	$Q_\pi(N)$	$Q_\pi(E)$	$Q_\pi(S)$	$Q_\pi(W)$	$Q_\pi(\text{exit})$
u	exit	NA	NA	NA	NA	1
w	W 	0	We converged on a local optimum!			1 NA
x	E 	-1			0 NA	
y	exit	N			NA	0
z	exit	NA	NA	NA	NA	-1

Must explore!

- Hard policy (insufficient): $\pi(s) = a$, $\pi: S \rightarrow \mathcal{A}$
- Soft policy: $\pi(a|s) = [0,1]$, $\pi: S \times \mathcal{A} \rightarrow p$
 - At the beginning $\forall a$, $\pi(a|s) > 0$ to allow exploration
 - Gradually shift towards a deterministic policy
- For instance: select a random action with probability ε
 - $\forall a \neq A^*$, $\pi(s, a) = \frac{\varepsilon}{|\mathcal{A}(s)|}$
 - Else select the greedy action: $\pi(s, A^*) = 1 - \varepsilon + \frac{\varepsilon}{|\mathcal{A}(s)|}$

ε -greedy MC control

On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow$ average($Returns(S_t, A_t)$)

$A^* \leftarrow \operatorname{argmax}_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

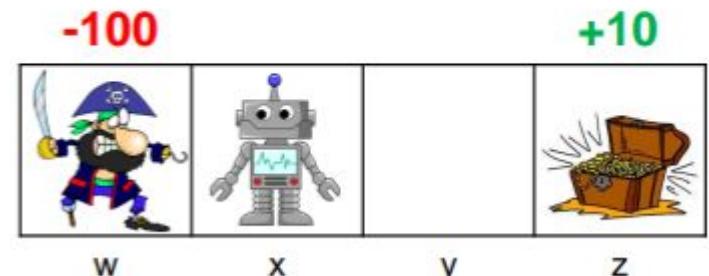
$$\gamma = 0.9$$

MC control - example

- $Q = \begin{array}{|c|c|c|c|} \hline 5 & 4,3 & 2,1 & 0 \\ \hline w & x & y & z \\ \hline \end{array}$

- $Returns = \begin{array}{|c|c|c|c|} \hline - & -, - & -, - & - \\ \hline w & x & y & z \\ \hline \end{array}$

- $\pi(a|s) = (1 - \varepsilon) \cdot \begin{array}{|c|c|c|c|} \hline \text{exit} & \leftarrow & \leftarrow & \text{exit} \\ \hline w & x & y & z \\ \hline \end{array}$



On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$A^* \leftarrow \arg\max_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

$$\gamma = 0.9$$

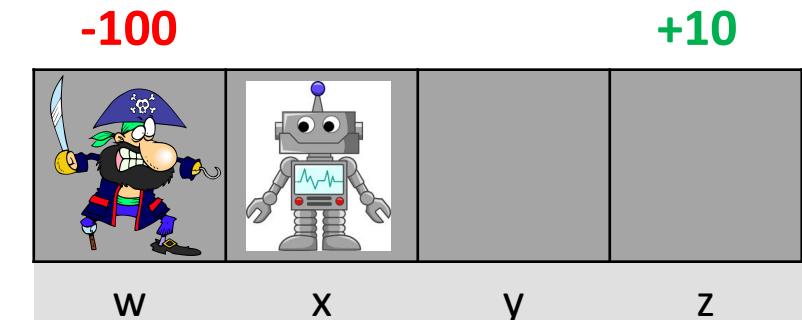
MC control - example

- $Q = \begin{array}{|c|c|c|c|} \hline 5 & 4,3 & 2,1 & 0 \\ \hline w & x & y & z \\ \hline \end{array}$

- $Returns = \begin{array}{|c|c|c|c|} \hline - & -, - & -, - & - \\ \hline w & x & y & z \\ \hline \end{array}$

- $\pi(a|s) = (1 - \varepsilon) \cdot \begin{array}{|c|c|c|c|} \hline \text{exit} & \rightarrow & \leftarrow & \text{exit} \\ \hline w & x & y & z \\ \hline \end{array}$

- $\tau = x, \leftarrow, 0, w, \text{exit}, -100$



On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$$G \leftarrow \gamma G + R_{t+1}$$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$$

$$A^* \leftarrow \underset{a}{\operatorname{argmax}} Q(S_t, a) \quad (\text{with ties broken arbitrarily})$$

For all $a \in \mathcal{A}(S_t)$:

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

$$\gamma = 0.9$$

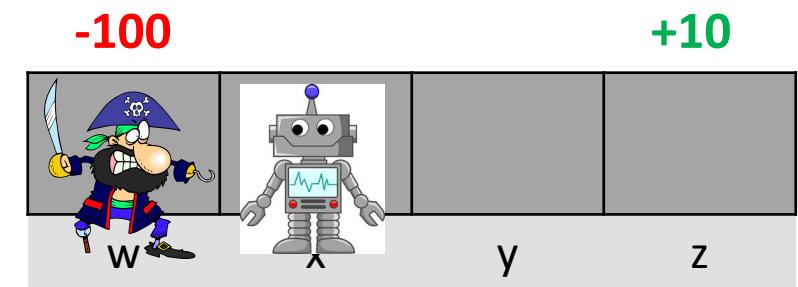
MC control - example

- $Q = \begin{array}{|c|c|c|c|} \hline 5 & 4,3 & 2,1 & 0 \\ \hline w & x & y & z \\ \hline \end{array}$

- $Returns = \begin{array}{|c|c|c|c|} \hline -100 & -90,0 & -, - & - \\ \hline w & x & y & z \\ \hline \end{array}$

- $\pi(a|s) = (1 - \varepsilon) \cdot \begin{array}{|c|c|c|c|} \hline \text{exit} & \rightarrow & \leftarrow & \text{exit} \\ \hline w & x & y & z \\ \hline \end{array}$
- $\varepsilon \cdot \text{Random}$

- $\tau = x, \leftarrow, 0, w, \text{exit}, -100$



On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$A^* \leftarrow \arg\max_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

$$\gamma = 0.9$$

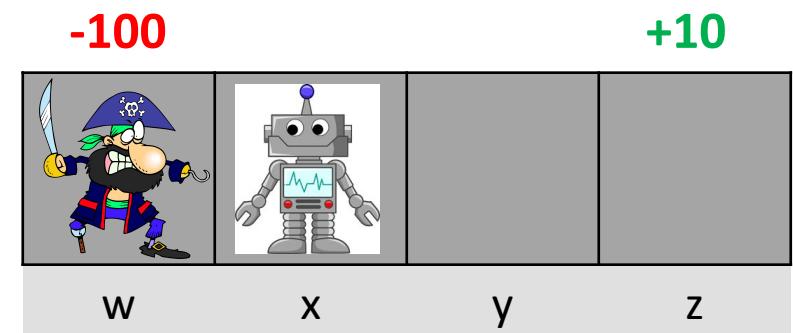
MC control - example

- $Q = \begin{array}{|c|c|c|c|} \hline -100 & -90,3 & 2,1 & 0 \\ \hline w & x & y & z \\ \hline \end{array}$

- $\overline{Returns} = \begin{array}{|c|c|c|c|} \hline -100 & -90,0 & -, - & - \\ \hline w & x & y & z \\ \hline \end{array}$

- $\pi(a|s) = (1 - \varepsilon) \cdot \begin{array}{|c|c|c|c|} \hline \text{exit} & \rightarrow & \leftarrow & \text{exit} \\ \hline w & x & y & z \\ \hline \end{array}$
- $\varepsilon \cdot \text{Random}$

- $\tau = x, \leftarrow, 0, w, \text{exit}, -100$



On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$A^* \leftarrow \arg\max_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$

$$\gamma = 0.9$$

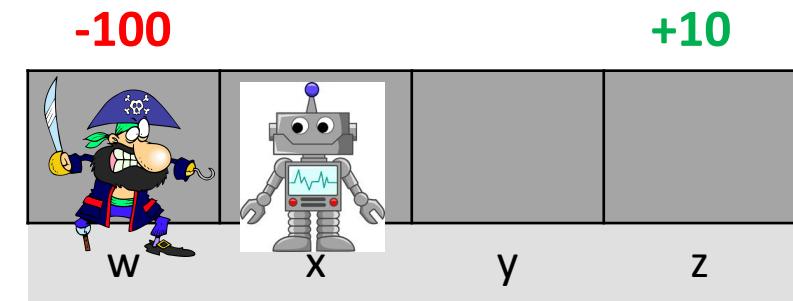
MC control - example

- $Q = \begin{array}{|c|c|c|c|} \hline -100 & -90,3 & 2,1 & 0 \\ \hline w & x & y & z \\ \hline \end{array}$

- $Returns = \begin{array}{|c|c|c|c|} \hline -100 & -90,0 & -, - & - \\ \hline w & x & y & z \\ \hline \end{array}$

- $\pi(a|s) = (1 - \varepsilon) \cdot \begin{array}{|c|c|c|c|} \hline \text{exit} & \rightarrow & \leftarrow & \text{exit} \\ \hline w & x & y & z \\ \hline \end{array}$
- $\varepsilon \cdot \text{Random}$

- $\tau = x, \leftarrow, 0, w, \text{exit}, -100$
- $A^* = [\rightarrow, \text{exit}]$



On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$A^* \leftarrow \arg\max_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

$$\gamma = 0.9$$

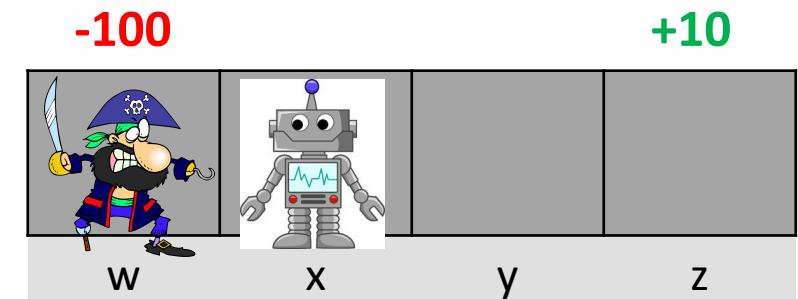
MC control - example

- $Q = \begin{array}{|c|c|c|c|} \hline -100 & -90,3 & 2,1 & 0 \\ \hline w & x & y & z \\ \hline \end{array}$

- $Returns = \begin{array}{|c|c|c|c|} \hline -100 & -90,0 & -, - & - \\ \hline w & x & y & z \\ \hline \end{array}$

- $\pi(a|s) = (1 - \varepsilon) \cdot \begin{array}{|c|c|c|c|} \hline \text{exit} & \rightarrow & \leftarrow & \text{exit} \\ \hline w & x & y & z \\ \hline \end{array}$
- $\varepsilon \cdot \text{Random}$

- $\tau = x, \leftarrow, 0, w, \text{exit}, -100$
- $A^* = [\rightarrow, \text{exit}]$



On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$A^* \leftarrow \text{argmax}_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

$$\gamma = 0.9$$

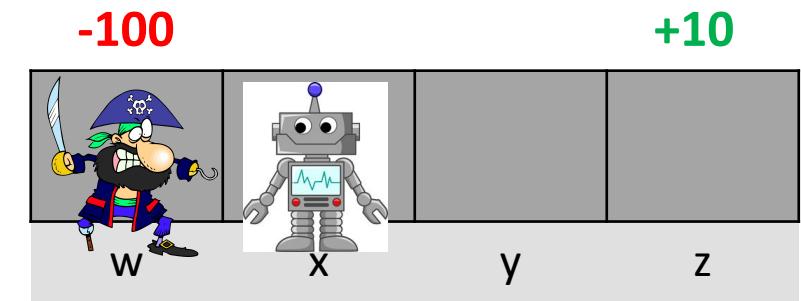
MC control - example

- $Q = \begin{array}{|c|c|c|c|} \hline -100 & -90,3 & 2,1 & 0 \\ \hline w & x & y & z \\ \hline \end{array}$

- $\overline{Returns} = \begin{array}{|c|c|c|c|} \hline -100 & -90,0 & -, - & - \\ \hline w & x & y & z \\ \hline \end{array}$

- $\pi(a|s) = (1 - \varepsilon) \cdot \begin{array}{|c|c|c|c|} \hline \text{exit} & \rightarrow & \leftarrow & \text{exit} \\ \hline w & x & y & z \\ \hline \end{array}$

- $\tau = x, \rightarrow, 0, y, \leftarrow, 0, x, \leftarrow, 0, \text{exit}, -100$



On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$A^* \leftarrow \text{argmax}_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

$$\gamma = 0.9$$

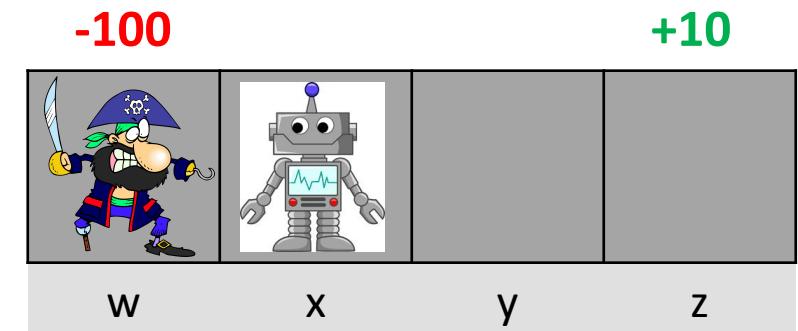
MC control - example

- $Q = \begin{array}{|c|c|c|c|} \hline -100 & -90, -72.9 & -81, 1 & 0 \\ \hline w & x & y & z \\ \hline \end{array}$

- $Returns = \begin{array}{|c|c|c|c|} \hline -100 & -90, -72.9 & -81, - & - \\ \hline w & x & y & z \\ \hline \end{array}$

- $\pi(a|s) = (1 - \varepsilon) \cdot \begin{array}{|c|c|c|c|} \hline \text{exit} & \rightarrow & \leftarrow & \text{exit} \\ \hline w & x & y & z \\ \hline \end{array}$
 - $\varepsilon \cdot \text{Random}$

- $\tau = x, \rightarrow, 0, y, \leftarrow, 0, x, \leftarrow, 0, \text{exit}, -100$



On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$A^* \leftarrow \arg \max_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$

$$\gamma = 0.9$$

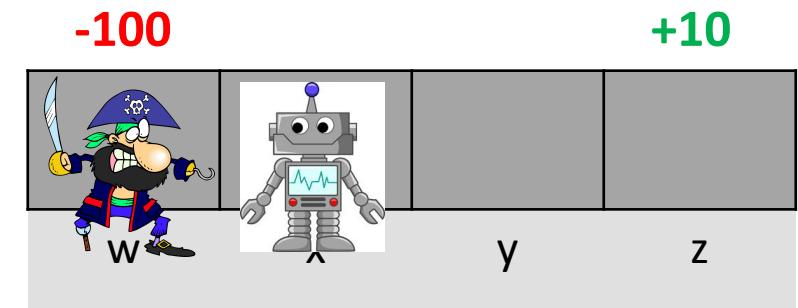
MC control - example

- $Q = \begin{array}{|c|c|c|c|} \hline -100 & -90, -72.9 & -81, 1 & 0 \\ \hline w & x & y & z \\ \hline \end{array}$

- $Returns = \begin{array}{|c|c|c|c|} \hline -100 & -90, -72.9 & -81, - & - \\ \hline w & x & y & z \\ \hline \end{array}$

- $\pi(a|s) = (1 - \varepsilon) \cdot \begin{array}{|c|c|c|c|} \hline \text{exit} & \rightarrow & \leftarrow & \text{exit} \\ \hline w & x & y & z \\ \hline \end{array}$
 - $\varepsilon \cdot \text{Random}$

- $\tau = x, \rightarrow, 0, y, \leftarrow, 0, x, \leftarrow, 0, \text{exit}, -100$
- $A^* = [\rightarrow, \rightarrow, \text{exit}]$



On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$A^* \leftarrow \arg \max_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon / |\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon / |\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

$$\gamma = 0.9$$

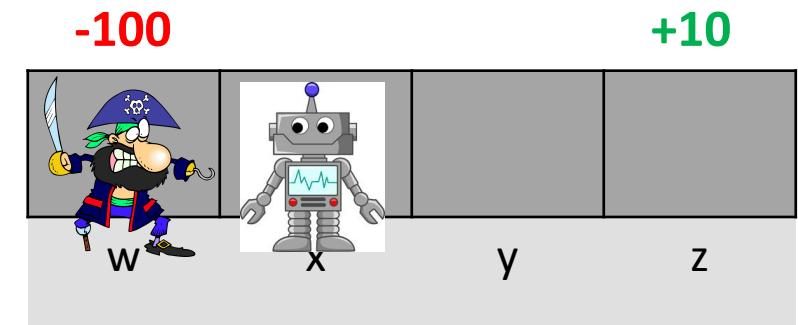
MC control - example

- $Q = \begin{array}{|c|c|c|c|} \hline -100 & -90, -72.9 & -81, 1 & 0 \\ \hline w & x & y & z \\ \hline \end{array}$

- $Returns = \begin{array}{|c|c|c|c|} \hline -100 & -90, -72.9 & -81, - & - \\ \hline w & x & y & z \\ \hline \end{array}$

- $\pi(a|s) = (1 - \varepsilon) \cdot \begin{array}{|c|c|c|c|} \hline \text{exit} & \rightarrow & \rightarrow & \text{exit} \\ \hline w & x & y & z \\ \hline \end{array}$
 - $\varepsilon \cdot \text{Random}$

- $\tau = x, \rightarrow, 0, y, \leftarrow, 0, x, \leftarrow, 0, \text{exit}, -100$
- $A^* = [\rightarrow, \rightarrow, \text{exit}]$



On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$A^* \leftarrow \arg\max_a Q(S_t, a)$ (with ties broken arbitrarily)

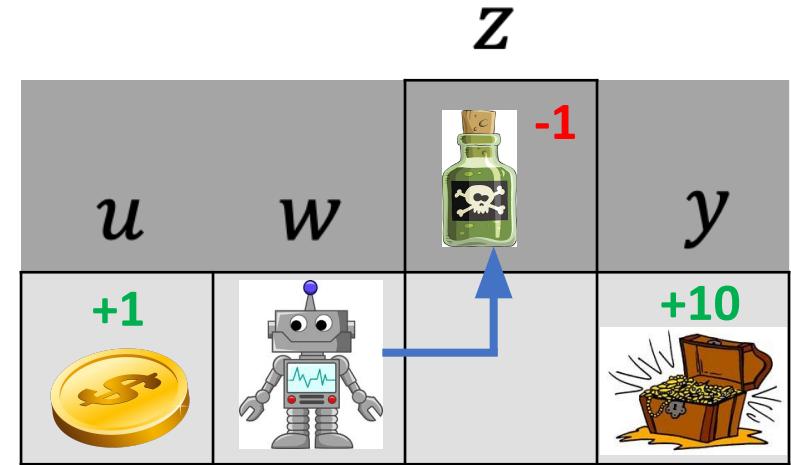
For all $a \in \mathcal{A}(S_t)$:

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

On-policy learning

Estimation True value

- $Q_\pi(w, \rightarrow) = q_\pi(w, \rightarrow) = -1$ (correct!)
- $q^*(w, \rightarrow) = ?$
- $q_{\pi \neq b} = ?$
- Observation drawn from π are useful for evaluating q_π
- Once the policy is changed these observations are irrelevant
- This is not **sample efficient!**



State	$\pi(s)$	$Q_\pi(\uparrow)$	$Q_\pi(\rightarrow)$
u	exit	NA	NA
w	\rightarrow	0	-1
x	\uparrow	-1	0
y	exit	NA	NA
z	exit	NA	NA

Quick Recap !

On-policy vs. Off-policy Learning

Off-policy learning

- We would like to use observations drawn from some policy b to evaluate q_π where $\pi \neq b$, specifically, we want to evaluate q_{π^*}
- We strive for full utilization of previous experience
- Off-policy learning allows us to optimize a ***target policy*** while following another ***behavior policy***
- **Pros:** sample efficient
- **Cons:** greater variance in value estimations

Off-policy learning conditions

- **Objective:** use episodes from b to estimate values for π
- For off-policy learning we must assume ***coverage***
 - $\forall s, a, \pi(a|s) > 0 \Rightarrow b(a|s) > 0$
- If this is true, then by running b repeatedly we will eventually discover all possible trajectories for π
- If coverage is violated for some s, a , then no inference is possible regarding that state-action value

Trajectory probability

- An agent following policy b sampled the following trajectory
 - $\tau = \{S_0, A_0, R_1, S_1, A_1, R_2, \dots, A_{T-1}, R_T, S_T\}$
- What is the probability of sampling this trajectory?
 - $\Pr\{\tau|b\} = b(A_0|S_0)p(S_1|S_0, A_0)b(A_1|S_1)p(S_2|S_1, A_1) \dots b(A_{T-1}|S_{T-1})p(S_T|S_{T-1}, A_{T-1})$
 - $= \prod_{k=0}^{T-1} b(A_k|S_k)p(S_{k+1}|S_k, A_k)$
- Assume MC control, how should $Q_b(S_t, A_t)$ be updated?
 - $G_t = \sum_{k=t+1}^T \gamma^{k-t+1} R_k$
 - Append G to $Returns(S_t, A_t)$ with weight $\Pr_t\{\tau|b\}$
 - $Q_b(S_t, A_t) = \text{Weighted_AVG}(Returns(S_t, A_t))$
- Computing the weighted average based on the sample probability would reduce variance (eliminate impact of noisy sampling)
- But the model, $p(S_{K+1}|S_K, A_K)$, is unknown! So can't compute $\Pr\{\tau|b\}$
- Can we say anything about $\Pr\{\tau|\pi \neq b\}$?

$$\mathbb{E}[X \sim p] = \sum_x p(x)x$$

Importance sampling

- Given a trajectory τ drawn by running b
 - We can **define** (not compute) the probability $\Pr\{\tau|b\}$
 - We can also **define** $\Pr\{\tau|\pi\}$
 - Define the ***importance sampling ratio*** as: $\rho_t = \frac{\Pr\{\tau_t|\pi\}}{\Pr\{\tau_t|b\}}$
 - Can we **compute** ρ without a model, $p(S_{K+1}|S_K, A_K)$?
 - $$\rho_t = \frac{\prod_{k=t}^{T-1} \pi(A_k|S_k) p(S_{K+1}|S_K, A_K)}{\prod_{k=t}^{T-1} b(A_k|S_k) p(S_{K+1}|S_K, A_K)} = \prod_{k=t}^{T-1} \frac{\pi(A_k|S_k)}{b(A_k|S_k)}$$
 YES!

Importance sampling

- **Fact:** $\mathbb{E}_{\tau \sim b}[G_t | S_t = s] = v_b(s)$
- Importance sampling allows us to compute an unbiased estimate of $v_\pi(s)$ by running b
- **Claim:** $\mathbb{E}_{\tau \sim b}[\rho_t G_t | S_t = s] = v_\pi(s)$
- We set $v_\pi(s)$ to be a weighted average of observed returns (weighted by the importance ratio)
- Assume visiting state s over M episodes using policy b
 - s is first visited during time step t^m during each episode, $m \in M$
- $v_\pi(s) = \frac{\sum_{m \in M} \rho_{t^m} G_{t^m}^m}{M}$

Importance sampling: proof

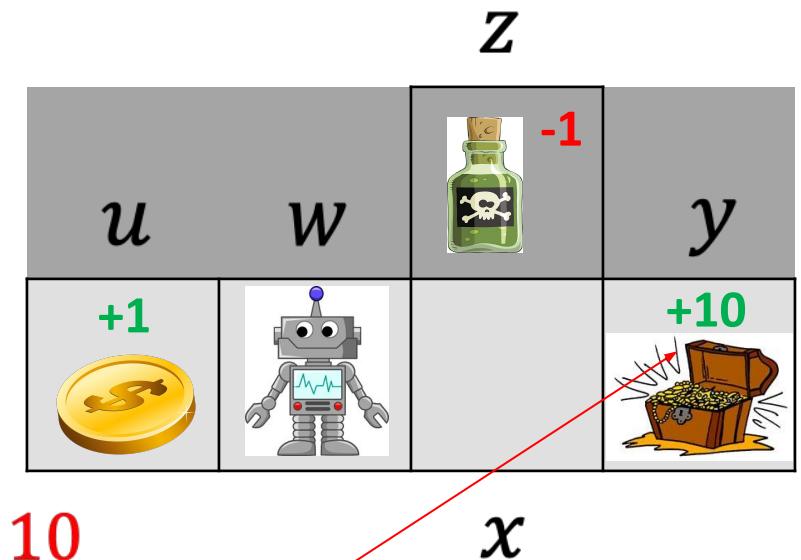
- We would like to estimate $\mathbb{E}[X]$ in our case when $X = v_\pi(s)$
- By definition: $\mathbb{E}[X] = \sum_x \Pr_X(x)x$
 - for the continuous case replace the sum with an integral
- If we don't know $\Pr(x)$ we can take a sample-based approach
- $\mathbb{E}[X] = \sum_x \Pr_X(x)x = \frac{\sum_{i=0}^n x_i}{n}$
 - This is an unbiased estimation because the samples are coming from the same distribution that defines \Pr_X
- But what if our samples come from a different distribution \Pr_Y ?

Importance sampling: proof

- Assume we know: $\mathbb{E}[Y] = \sum_y \Pr(y)y = \frac{\sum_{i=0}^n y_i}{n}$
 - Can we use this to compute $\mathbb{E}[X]$?
 - Yes! if we assume that all possible values of X exist in Y (*converge*)
 - $\mathbb{E}[X] = \sum_x \Pr_X(x)x = \sum_y \Pr_X(y)y$
 - $= \sum_y \frac{\Pr_X(y)}{\Pr_Y(y)} \Pr_Y(y)y$
 - $= \frac{\Pr_X(y)}{\Pr_Y(y)} \cdot \frac{\sum_{i=0}^n y_i}{n}$
- Importance ratio

(ordinary) Importance sampling - example

- $\because b(s|a) = \begin{cases} \rightarrow, p(0.5) \\ \leftarrow, p(0.5) \end{cases}$
- $\bullet \pi(s|a) = \begin{cases} \rightarrow, p(0.99) \\ \leftarrow, p(0.01) \end{cases}$
- $\bullet \tau_1 = \{w, \rightarrow, 0, x, \rightarrow, 0, y, exit, 10\}$
- $\bullet v_\pi(w) = \frac{\sum_{m \in M} \rho_t^m G_t^m}{M} = \frac{0.99}{0.5} * \frac{0.99}{0.5} * 10 = \textcolor{red}{3.96 * 10}$
- $\bullet \tau_2 = \{w, \leftarrow, 0, u, exit, 1\}$
- $\bullet v_\pi(w) = \frac{\sum_{m \in M} \rho_t^m G_t^m}{M} = \frac{3.96*10+0.02*1}{2}$

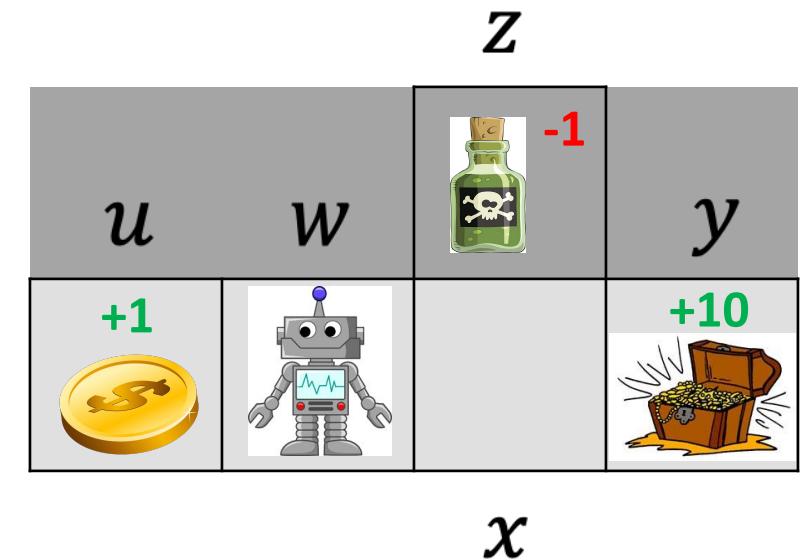


39.6 ??
 Ordinary Importance
 sampling is unbiased yet
 high variance

Weighted importance sampling

- $v_\pi(s) = \frac{\sum_{m \in M} [\rho_t^m G_t^m]}{\sum_{m \in M} \rho_t^m}$
- $b(s|a) = \begin{cases} \rightarrow, p(0.5) \\ \leftarrow, p(0.5) \end{cases}$
- $\pi(s|a) = \begin{cases} \rightarrow, p(0.99) \\ \leftarrow, p(0.01) \end{cases}$
- $\tau_1 = \{w, \rightarrow, 0, x, \rightarrow, 0, y, exit, 10\}$
- $v_\pi(w) = \frac{\sum_{m \in M} [\rho_t^m G_t^m]}{\sum_{m \in M} \rho_t^m} = \frac{3.96*10}{3.96}$
- $\tau_2 = \{w, \leftarrow, 0, u, exit, 1\}$
- $v_\pi(w) = \frac{\sum_{m \in M} [\rho_t^m G_t^m]}{\sum_{m \in M} \rho_t^m} = \frac{3.96*10 + 0.02*1}{3.98}$

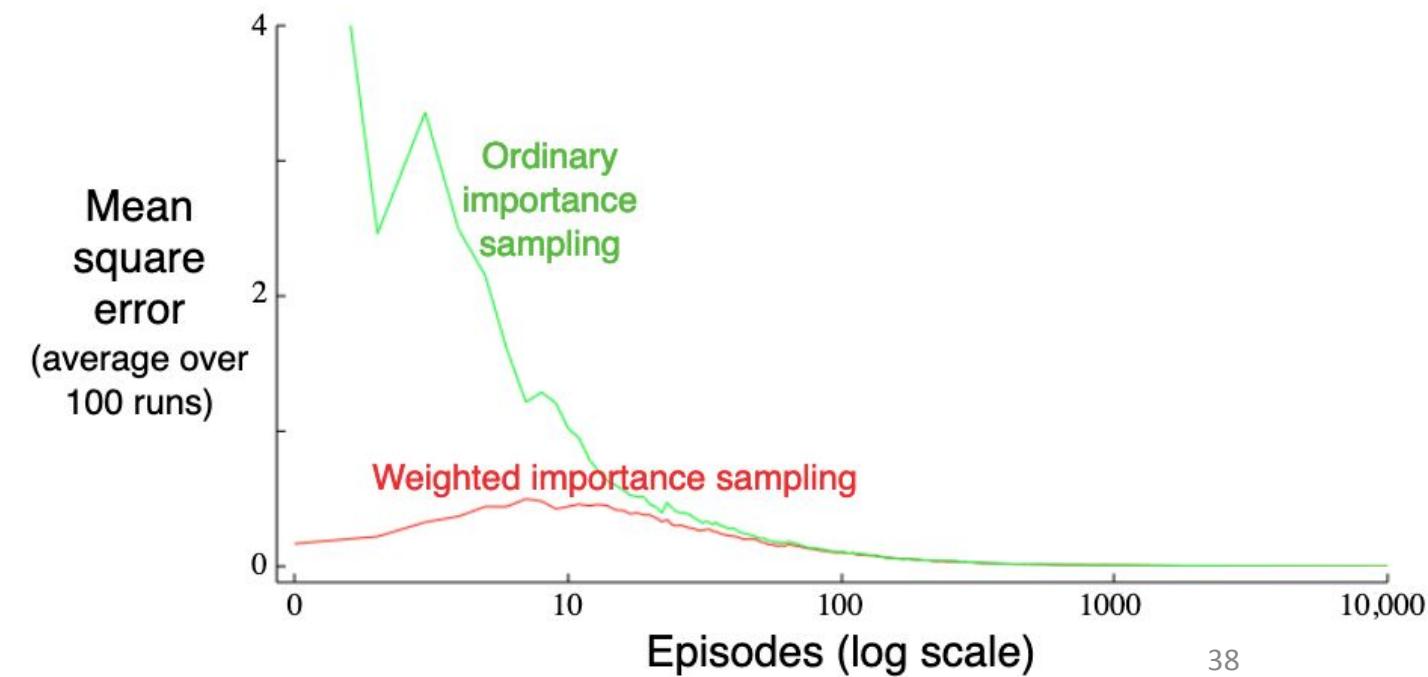
Trick: normalize by the sum of importance ratios



Ordinary Importance sampling is unbiased while the weighted version is biased (initially). Ordinary Importance sampling results in high variance while the weighted version has a bounded variance

Ordinary vs weighted importance sampling

- Estimating a black-jack state
- Target policy: hit on 19 or below
- Behavior policy: random (uniform)
- Both approaches converge to the true value
- weighted importance sampling is much better initially



Revisit : Incremental Implementation (Slide taken from Multi-armed Bandit)

Note:

- StepSize decreases with each update
- We use α or $\alpha_t(a)$ to denote step size (constant / varies with each step)

Discussion:

Const vs. Variable step size?

$$\begin{aligned}
 Q_{n+1} &= \frac{1}{n} \sum_{i=1}^n R_i \\
 &= \frac{1}{n} \left(R_n + \sum_{i=1}^{n-1} R_i \right) \\
 &= \frac{1}{n} \left(R_n + (n-1) \frac{1}{n-1} \sum_{i=1}^{n-1} R_i \right) \\
 &= \frac{1}{n} \left(R_n + (n-1)Q_n \right) \\
 &= \frac{1}{n} \left(R_n + nQ_n - Q_n \right) \\
 &= Q_n + \frac{1}{n} [R_n - Q_n],
 \end{aligned}$$

$$\text{NewEstimate} \leftarrow \text{OldEstimate} + \text{StepSize} [\text{Target} - \text{OldEstimate}]$$

MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T - 1, T - 2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

Accumulated
reward after step t

MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$$b \leftarrow \text{any soft policy}$$

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T - 1, T - 2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

1 over Joint probability over observed actions from following b after step t . This equals the IS ratio here because the target policy is deterministic.

MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T - 1, T - 2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

Going back in time

MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T - 1, T - 2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

Discount future rewards and add immediate reward

MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T - 1, T - 2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

Cumulative sum of IS weights affiliated with S_t, A_t (for weighted IS)

MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T - 1, T - 2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

Incremental update of Q values (waited moving average)

MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T - 1, T - 2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

Update target policy
(greedy)

MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$$b \leftarrow \text{any soft policy}$$

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T - 1, T - 2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

Since π is deterministic,
once we diverge from it all
IS weights of earlier
actions will be 0

MC control + importance sampling

Off-policy MC control, for estimating $\pi \approx \pi_*$

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$$b \leftarrow \text{any soft policy}$$

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T - 1, T - 2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

Update the joint prob by multiplying by ρ_t . Notice that $\pi(S_t) = 1$ in this example (deterministic target policy)

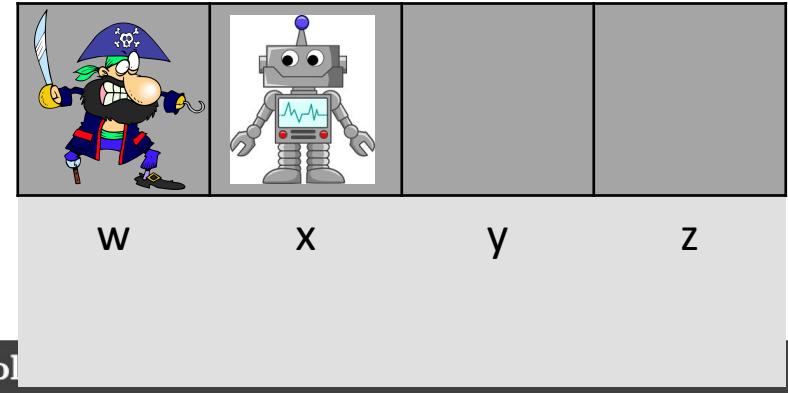
MC control + IS example

-

-	4,3	2,1	-
W	X	Y	Z
-	0,0	0,0	-
W	X	Y	Z
exit			exit
W	X	Y	Z

-100

+10



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

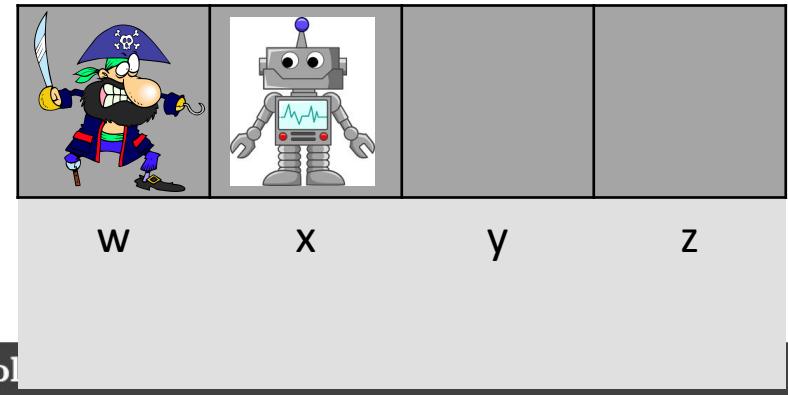
$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-	4,3	2,1	-
W	X	Y	Z
-	0,0	0,0	-
W	X	Y	Z
exit			exit
W	X	Y	Z

-100 +10



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

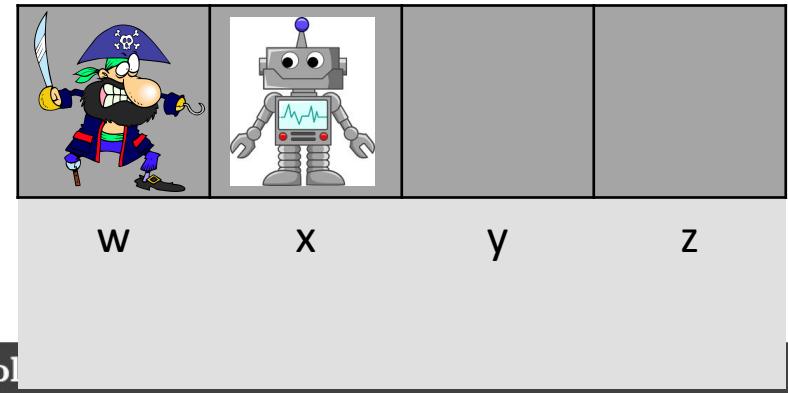
$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

5	4,3	2,1	5
w	x	y	z
0	0,0	0,0	0
w	x	y	z
exit			exit
w	x	y	z

-100 +10



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

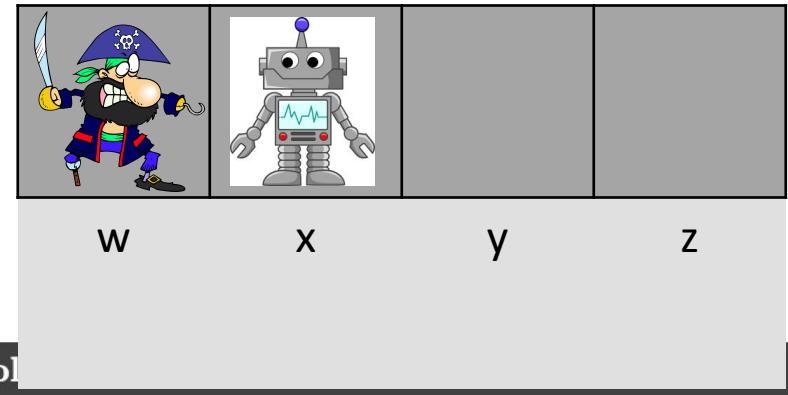
$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

5	4,3	2,1	5
w	x	y	z
0	0,0	0,0	0
w	x	y	z
exit			exit
w	x	y	z

-100 +10



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

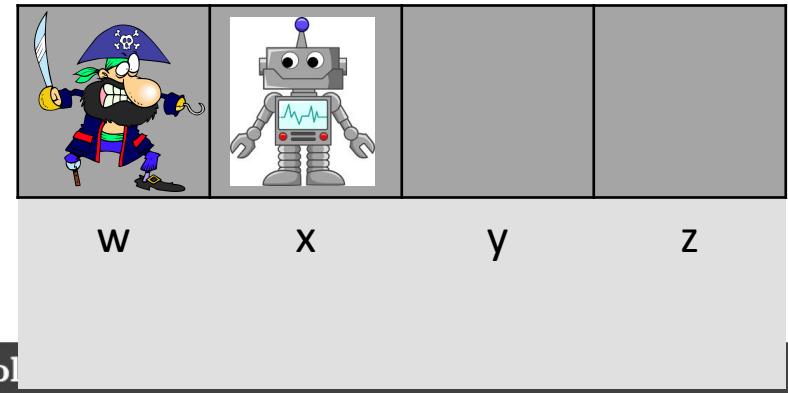
$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

5	4,3	2,1	5
w	x	y	z
1	0,0	0,0	0
w	x	y	z
exit			exit
w	x	y	z

-100 +10



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

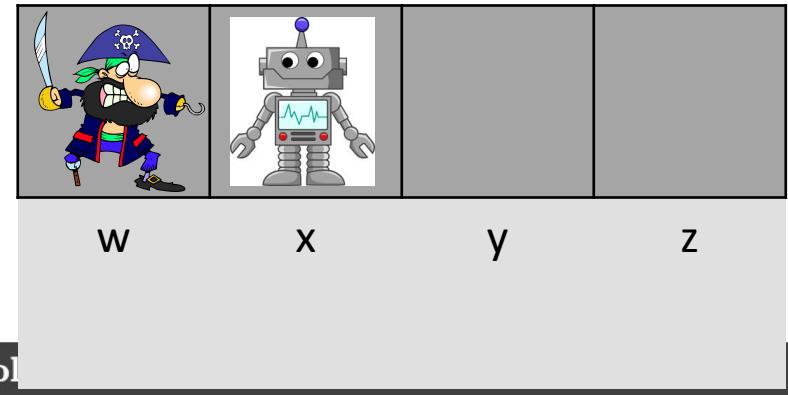
$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	4,3	2,1	5
w	x	y	z
1	0,0	0,0	0
w	x	y	z
exit			exit
w	x	y	z

-100 +10



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

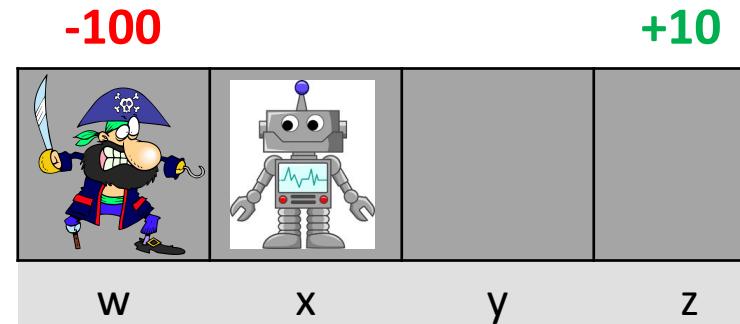
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	4,3	2,1	5
w	x	y	z
1	0,0	0,0	0
w	x	y	z
exit			exit
w	x	y	z



Off-policy

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T - 1, T - 2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

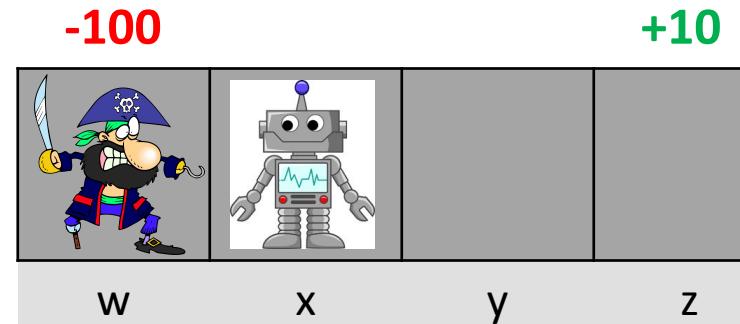
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	4,3	2,1	5
w	x	y	z
1	0,0	0,0	0
w	x	y	z
exit			exit
w	x	y	z



Off-policy

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

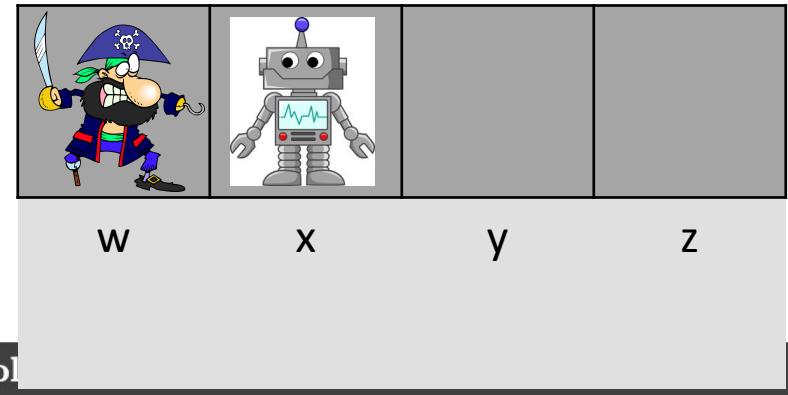
$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	4,3	2,1	5
w	x	y	z
1	0,0	0,0	0
w	x	y	z
exit			exit
w	x	y	z

-100 +10



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

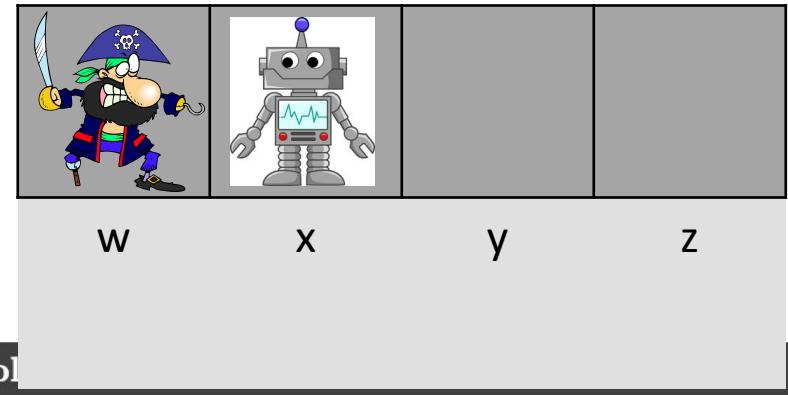
$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	4,3	2,1	5
w	x	y	z
1	2,0	0,0	0
w	x	y	z
exit			exit
w	x	y	z

-100 +10



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

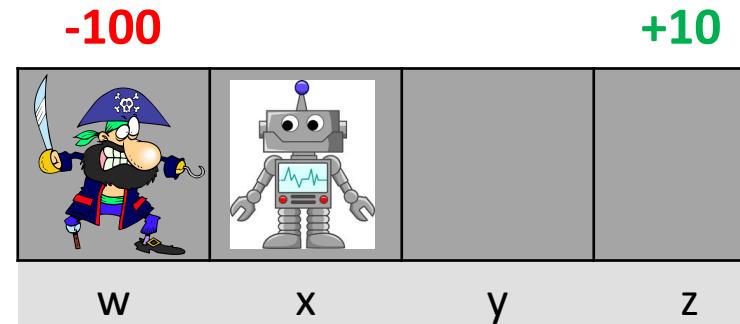
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,1	5
w	x	y	z
1	2,0	0,0	0
w	x	y	z
exit			exit
w	x	y	z



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

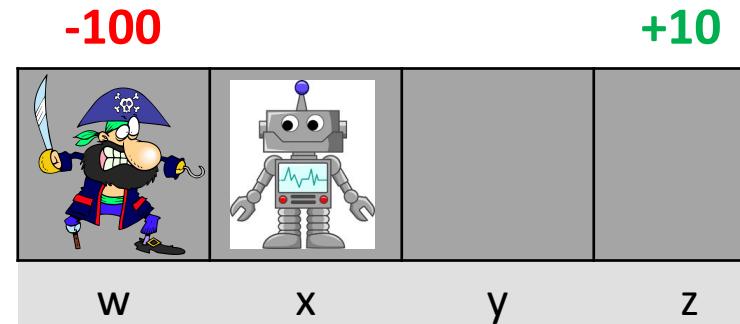
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,1	5
w	x	y	z
1	2,0	0,0	0
w	x	y	z
exit			exit
w	x	y	z



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

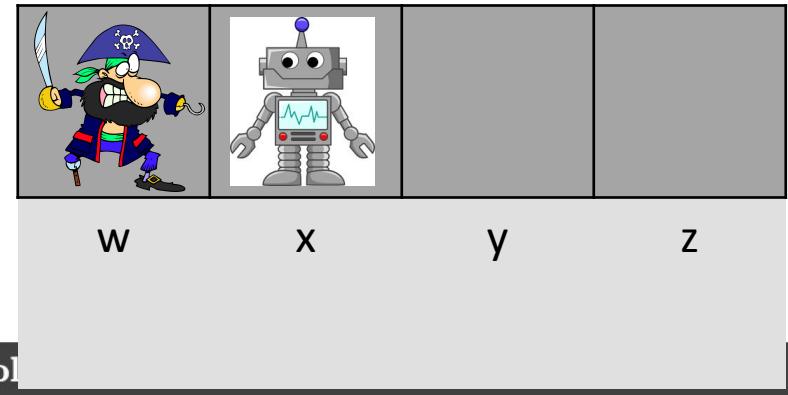
$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,1	5
w	x	y	z
1	2,0	0,0	0
w	x	y	z
exit			exit
w	x	y	z

-100 +10



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

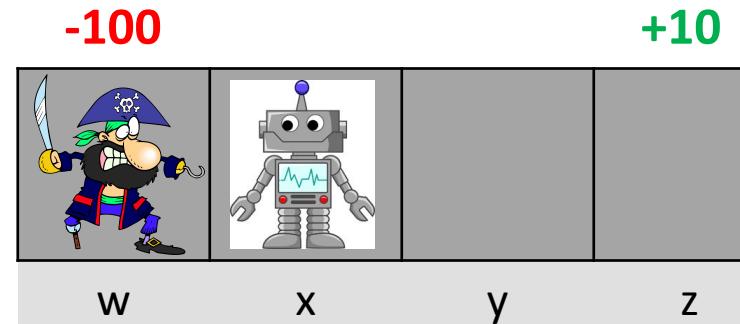
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,1	5
w	x	y	z
1	2,0	0,0	0
w	x	y	z
exit			exit
w	x	y	z



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

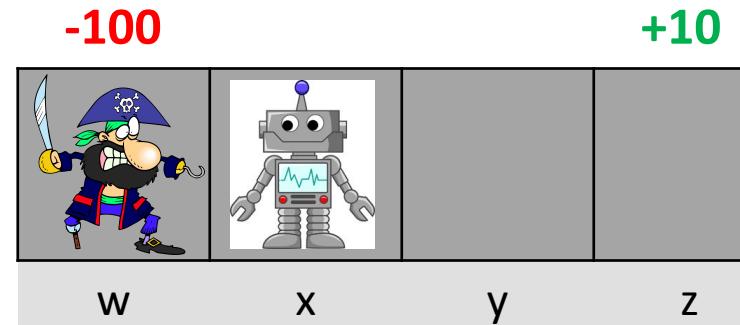
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,1	5
w	x	y	z
1	2,0	0,0	0
w	x	y	z
exit			exit
w	x	y	z



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

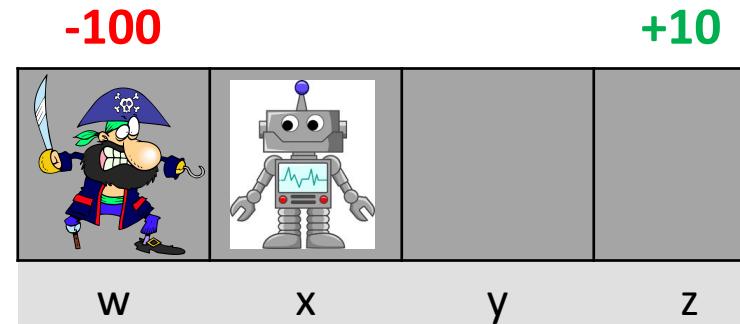
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,1	5
w	x	y	z
1	2,0	0,0	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

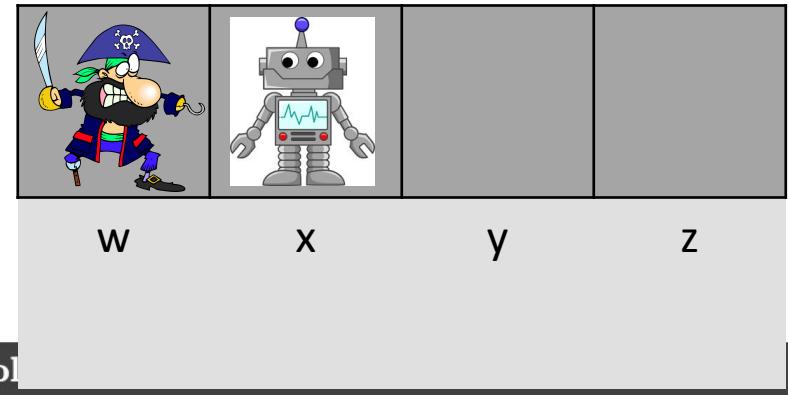
$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,1	10
w	x	y	z
1	2,0	0,0	1
w	x	y	z
exit			exit
w	x	y	z

-100 +10



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

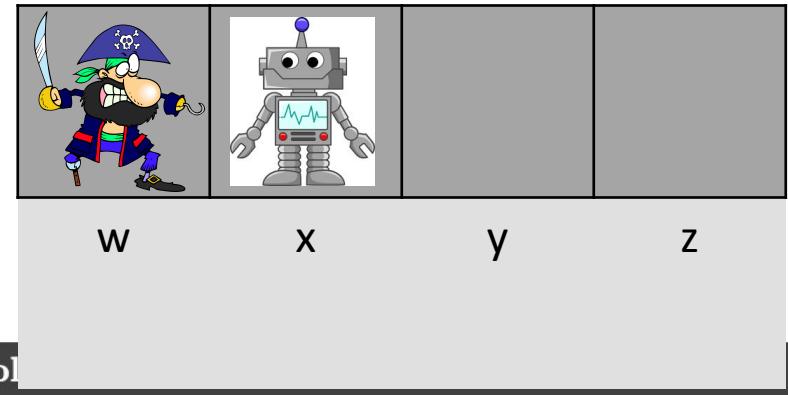
$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,1	10
w	x	y	z
1	2,0	0,0	1
w	x	y	z
exit			exit
w	x	y	z

-100 +10



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

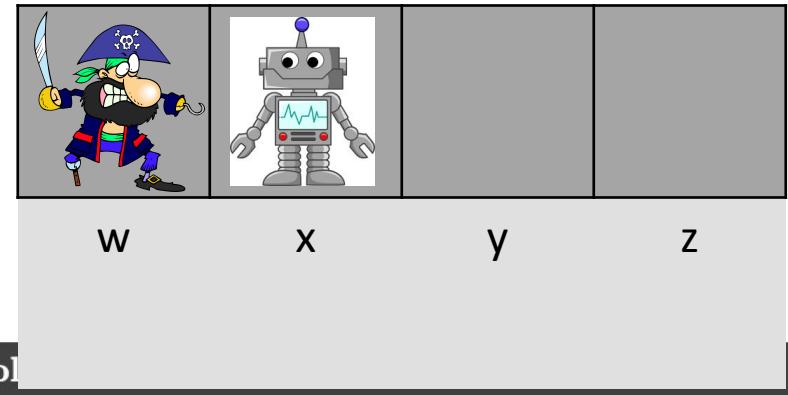
$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,1	10
w	x	y	z
1	2,0	0,0	1
w	x	y	z
exit			exit
w	x	y	z

-100 +10



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

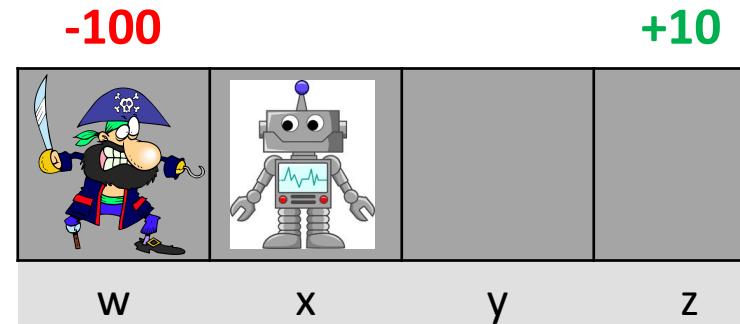
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,1	10
w	x	y	z
1	2,0	0,0	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T - 1, T - 2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

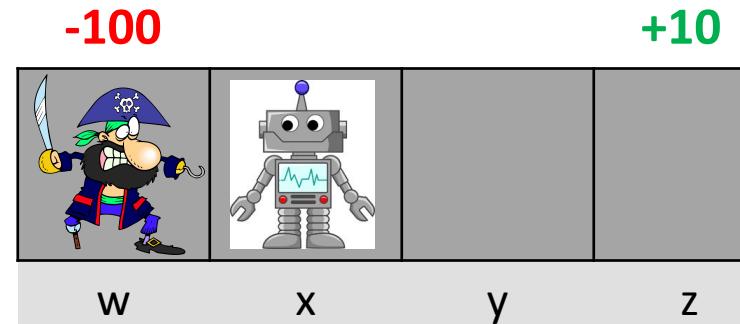
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,1	10
w	x	y	z
1	2,0	0,2	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

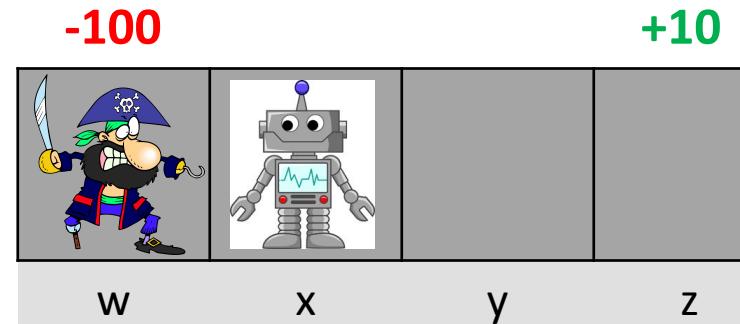
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,9	10
w	x	y	z
1	2,0	0,2	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

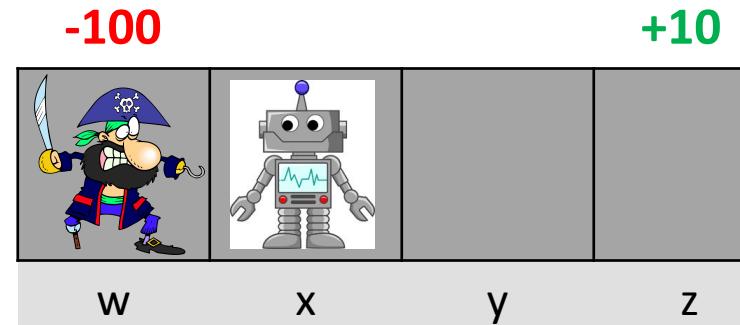
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,9	10
w	x	y	z
1	2,0	0,2	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

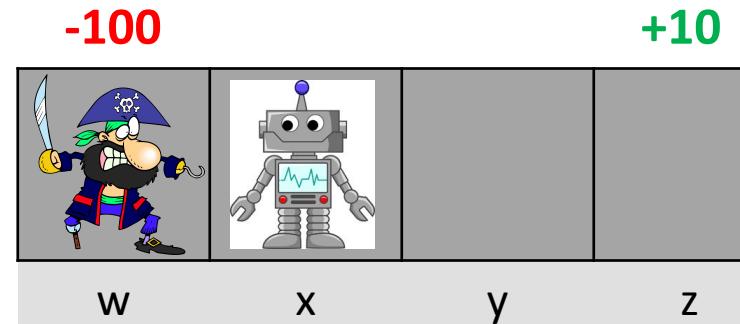
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,9	10
w	x	y	z
1	2,0	0,2	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

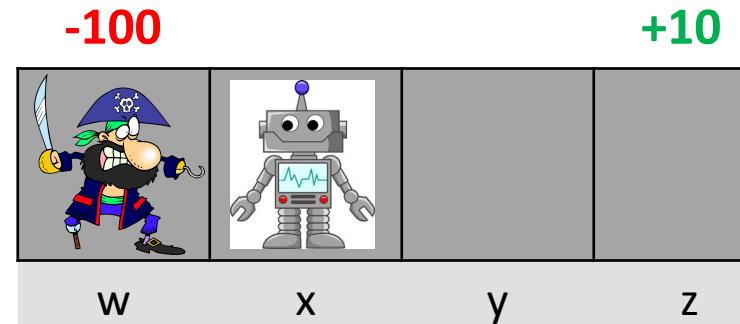
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,9	10
w	x	y	z
1	2,0	0,2	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

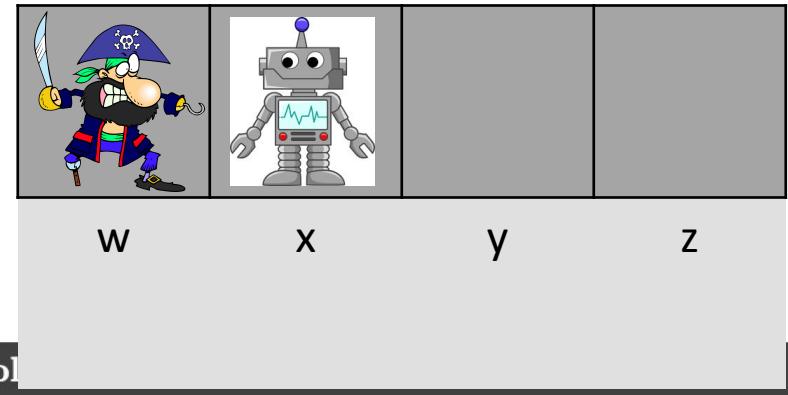
$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,3	2,9	10
w	x	y	z
1	2,4	0,2	1
w	x	y	z
exit			exit
w	x	y	z

-100 +10



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

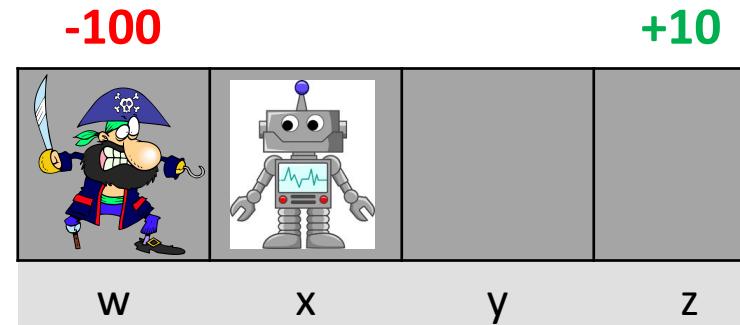
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,8.1	2,9	10
w	x	y	z
1	2,4	0,2	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

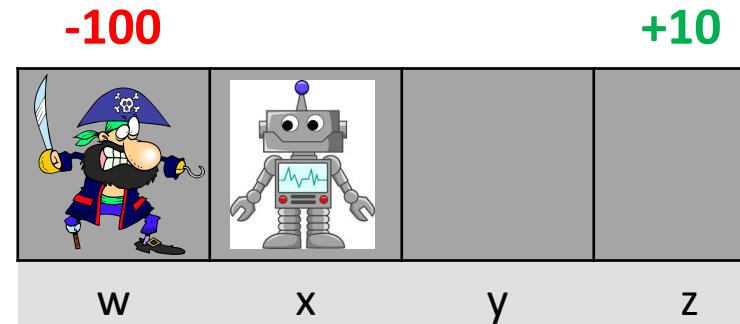
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,8.1	2,9	10
w	x	y	z
1	2,4	0,2	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

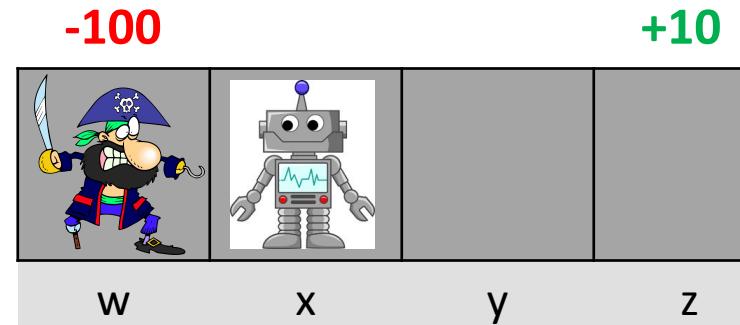
If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

MC control + IS example

-

-100	-90,8.1	2,9	10
w	x	y	z
1	2,4	0,2	1
w	x	y	z
exit			exit
w	x	y	z



Off-policy MC Control

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \leftarrow \text{arbitrary}$$

$$C(s, a) \leftarrow 0$$

$$\pi(s) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

Repeat forever:

$b \leftarrow$ any soft policy

Generate an episode using b :

$$S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T, S_T$$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For $t = T-1, T-2, \dots$ down to 0:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit For loop

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

What did we learn?

- Online evaluation through the Monte-Carlo approach
 - Update V or Q estimations based on observed returns
- On policy Monte-Carlo control
 - Beware of local optimum. Must explore!
 - Consider using a soft policy $\pi(a|s)$
 - Assign soft policy values that mimic an ε -greedy strategy
 - On policy learning is not sample efficient!
- Off policy Monte-Carlo control
 - Define target policy (e.g., $\arg \max_a Q(S_t, a)$) that can be different than the behavior policy
 - Utilize weighted importance sampling to train a target policy
 - $v_\pi(s) = \frac{\sum_{m \in M} [\rho_{tm} G_t^m]}{\sum_{m \in M} \rho_{tm}}$



Required Readings

1. Chapter-3,4 of Introduction to Reinforcement Learning, 2nd Ed., Sutton & Barto



BITS Pilani
Pilani | Dubai | Goa | Hyderabad

Thank you