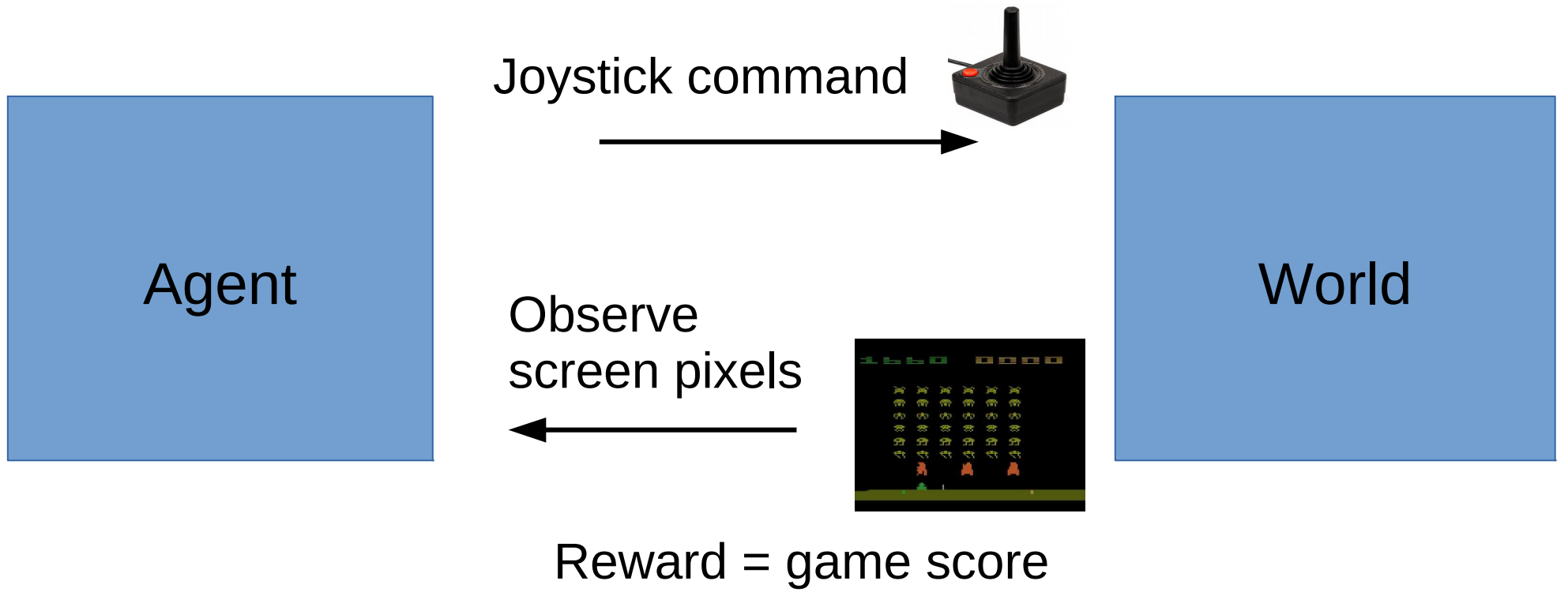# Monte Carlo Approaches to Reinforcement Learning

Robert Platt (w/ Marcus Gualtieri's edits)
Northeastern University
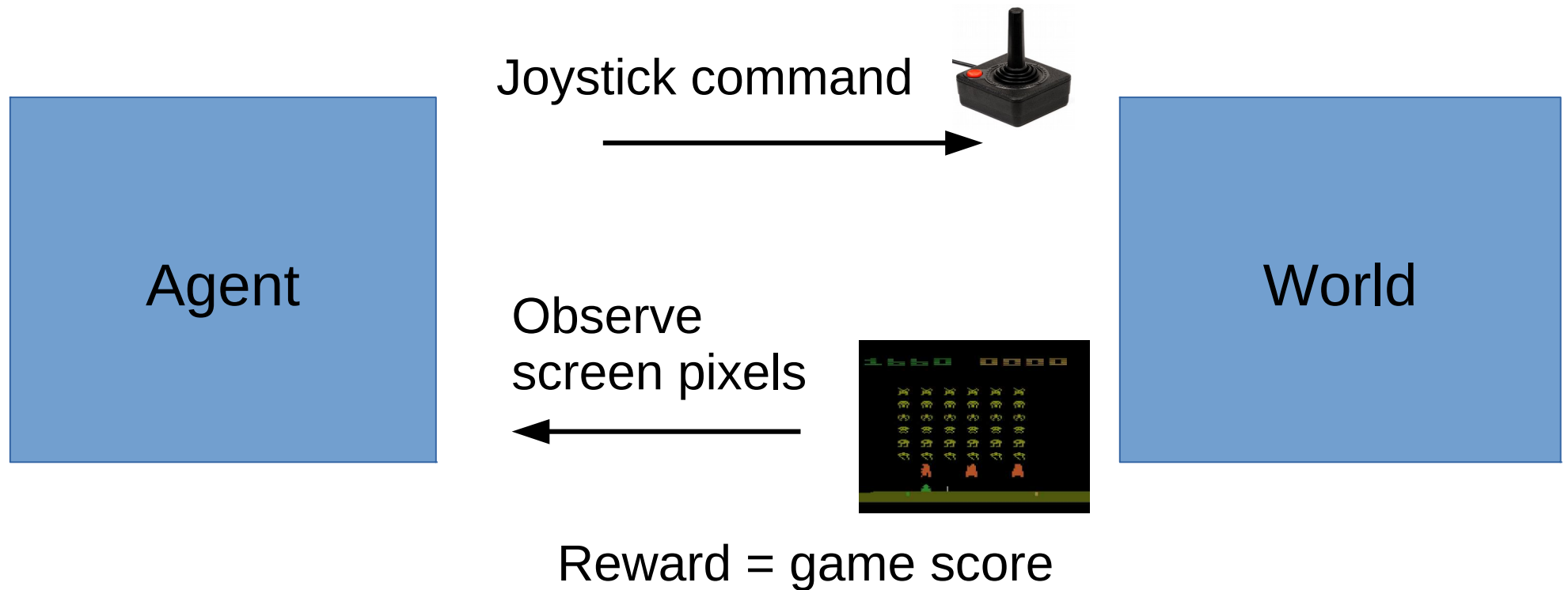
# Model Free Reinforcement Learning

Joystick command

Agent

World

Observe
screen pixels

Reward = game score

Goal: learn a value function through trial-and-error experience

# Model Free Reinforcement Learning

Joystick command

Agent

World

Observe
screen pixels

Reward = game score

Goal: learn a value function through trial-and-error experience

Recall: $V^\pi(s_t = s) \equiv$ Value of state $s$ when acting according to policy $\pi$

# Model Free Reinforcement Learning
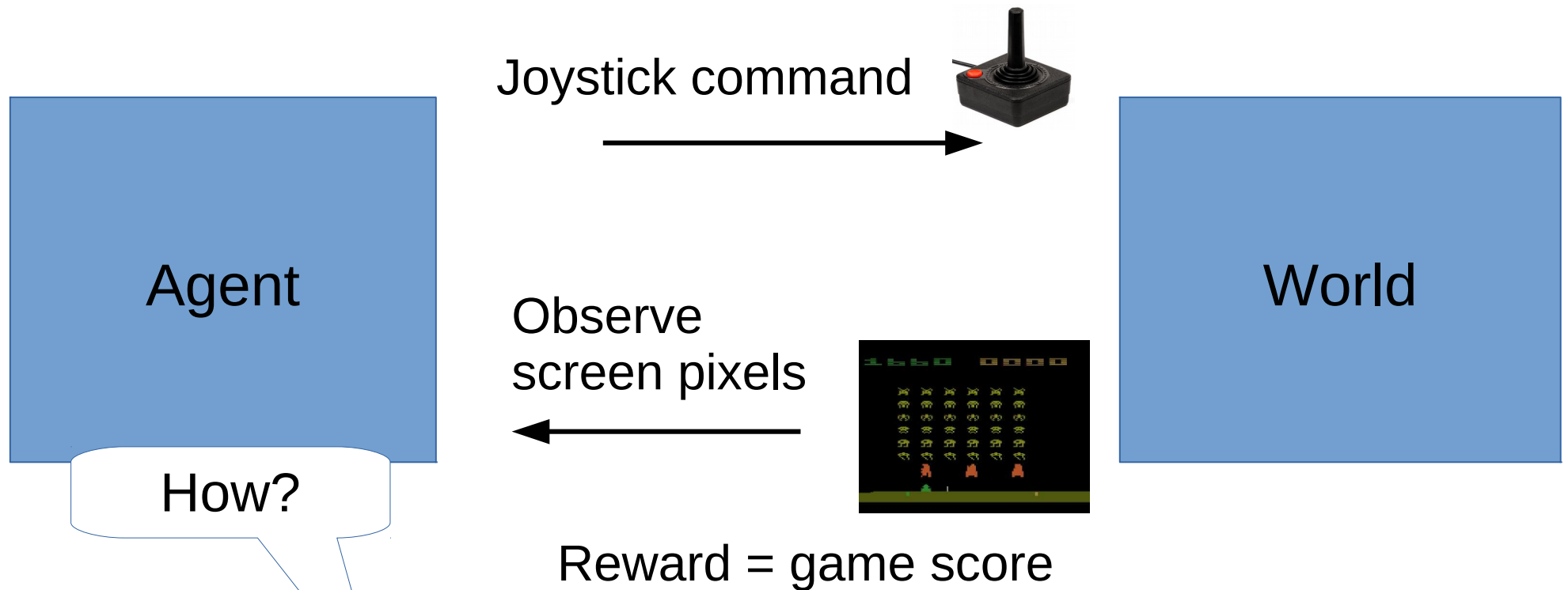


Joystick command

Agent

World

Observe
screen pixels

How?

Reward = game score

Goal: learn a value function through trial-and-error experience

Recall: $V^{\pi}(s_t = s) \equiv$ Value of state $s$ when acting according to policy $\pi$

# Model Free Reinforcement Learning

Agent

World

Simplest solution: average all outcomes
from previous experiences in a given state

– this is called a *Monte Carlo* method

...een pixels

How?

Reward = game score

Goal: learn a value function through trial-and-error experience

Recall: $V^\pi(s_t = s) \equiv$ Value of state $s$ when acting according to policy $\pi$

# Running Example: Blackjack





**State:** sum of cards in agent's hand + dealer's showing card + does agent have usable ace?

**Actions:** hit, stick

**Objective:** Have agent's card sum be greater than the dealer's without exceeding 21

**Reward:** +1 for winning, 0 for a draw, -1 for losing

**Discounting:** $\gamma = 1$

**Dealer policy:** draw until sum at least 17

# Running Example: Blackjack



Blackjack "Basic Strategy" is a set of rules for play so as to maximize return
– well known in the gambling community
– how might an RL agent *learn* the Basic Strategy?

# Monte Carlo Policy Evaluation: Example

Dealer card:

Agent's hand:

State         Action         Next State         Reward

# Monte Carlo Policy Evaluation: Example

Dealer card:

Agent's hand:

| State | Action | Next State | Reward |
|-------|--------|------------|--------|
| 19, 10, no | | | |

# Monte Carlo Policy Evaluation: Example

Dealer card:

Agent's hand:

Agent sum, dealer's card, ace?

State          Action          Next State          Reward

19, 10, no

# Monte Carlo Policy Evaluation: Example

Dealer card:

Agent's hand:

Agent sum, dealer's card, ace?

| State | Action | Next State | Reward |
| --- | --- | --- | --- |
| 19, 10, no | HIT | | |

# Monte Carlo Policy Evaluation: Example

Dealer card:



Agent's hand:



Agent sum, dealer's card, ace?

| State | Action | Next State | Reward |
|-------|--------|------------|--------|
| 19, 10, no | HIT | 22, 10, no | -1 |

# Monte Carlo Policy Evaluation: Example

Dealer card:

Agent's hand:

Agent sum, dealer's card, ace?

Bust!
(reward = -1)

| State | Action | Next State | Reward |
|-------|--------|------------|--------|
| 19, 10, no | HIT | 22, 10, no | -1 |

# Monte Carlo Policy Evaluation: Example

| State | Action | Next State | Reward |
|-------|--------|------------|--------|
| 19, 10, no | HIT | 22, 10, no | -1 |

Upon episode termination, make the following value function updates:

$$V((19, 10, no)) \leftarrow -1$$

# Monte Carlo Policy Evaluation: Example

Next episode...

# Monte Carlo Policy Evaluation: Example

Dealer card:

Agent's hand:

| State | Action | Next State | Reward |
|-------|--------|------------|--------|
| 13, 10, no | | | |

# Monte Carlo Policy Evaluation: Example

Dealer card:

Agent's hand:

| State | Action | Next State | Reward |
|---|---|---|---|
| 13, 10, no | HIT | 16, 10, no | 0 |

# Monte Carlo Policy Evaluation: Example

Dealer card:

Agent's hand:

| State | Action | Next State | Reward |
|-------|--------|------------|--------|
| 13, 10, no | HIT | 16, 10, no | 0 |
| 13, 10, no | | | |

# Monte Carlo Policy Evaluation: Example

Dealer card:

Agent's hand:

| State | Action | Next State | Reward |
|-------|--------|------------|--------|
| 13, 10, no | HIT | 16, 10, no | 0 |
| 13, 10, no | HIT | 19, 10, no | 0 |

# Monte Carlo Policy Evaluation: Example

Dealer card:



Agent's hand:



| State | Action | Next State | Reward |
|-------|--------|-----------|--------|
| 13, 10, no | HIT | 16, 10, no | 0 |
| 13, 10, no | HIT | 19, 10, no | 0 |
| 19, 10, no | | | |

# Monte Carlo Policy Evaluation: Example



Dealer card:

Agent's hand:

| State | Action | Next State | Reward |
|-------|--------|------------|--------|
| 13, 10, no | HIT | 16, 10, no | 0 |
| 13, 10, no | HIT | 19, 10, no | 0 |
| 19, 10, no | HIT | 21, 22, no | 1 |

# Monte Carlo Policy Evaluation: Example

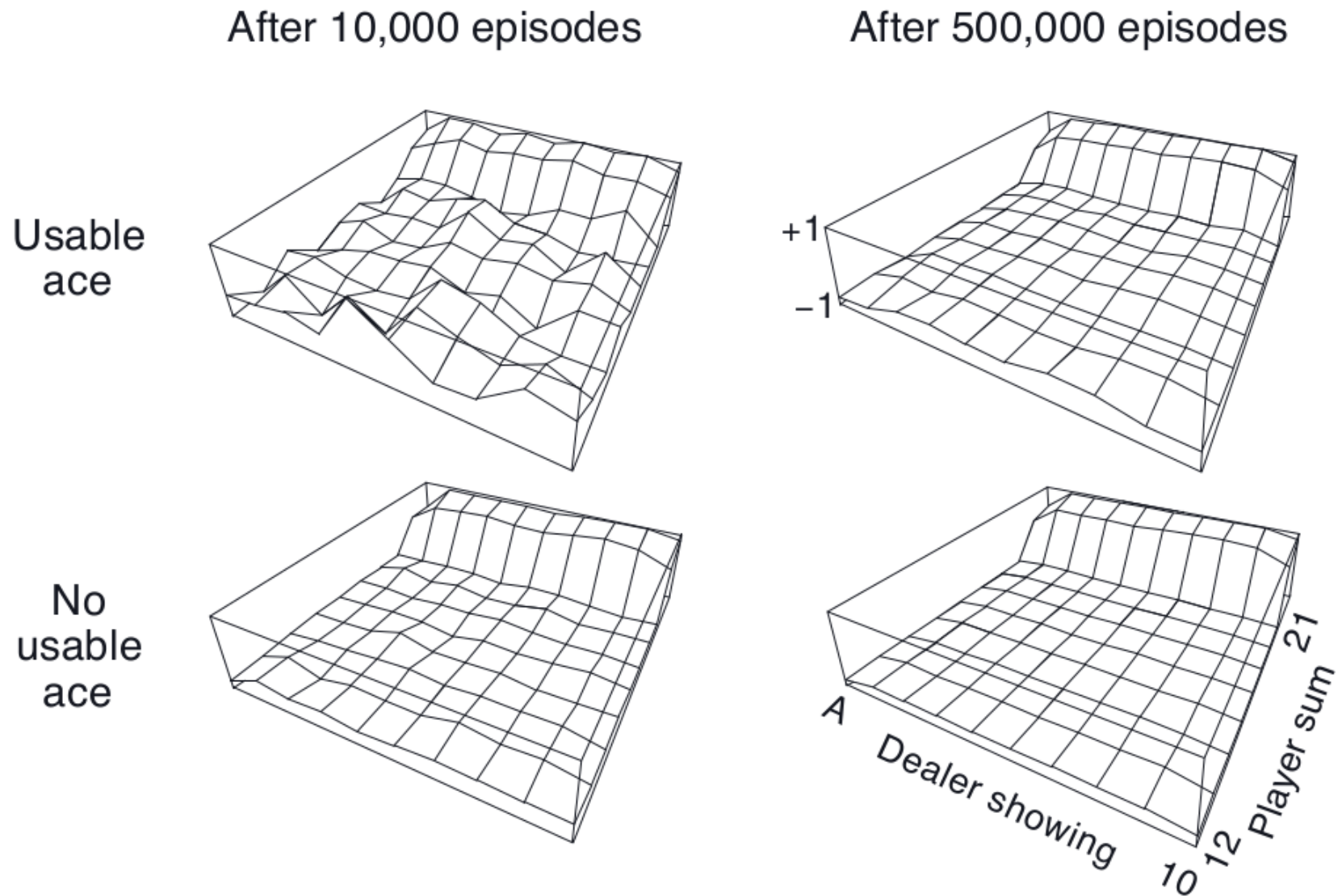| State | Action | Next State | Reward |
|-------|--------|------------|--------|
| 13, 10, no | HIT | 16, 10, no | 0 |
| 16, 10, no | HIT | 19, 10, no | 0 |
| 19, 10, no | HIT | 21, 22, no | 1 |

Upon episode termination, make the following value function updates:

$$V((13, 10, no)) \leftarrow 1$$

$$V((16, 10, no)) \leftarrow 1$$

$$V((19, 10, no)) \leftarrow (-1 + 1)/2$$

# Monte Carlo Policy Evaluation: Example



After 10,000 episodes     After 500,000 episodes

Usable ace

No usable ace

+1
−1

A

Dealer showing

10   12   21

Player sum

Value function learned for "hit everything except for 20 and 21" policy.

# Monte Carlo Policy Evaluation

Given a policy, $\pi$, estimate the value function, $V(s)$, for all states, $s \in \mathcal{S}$

# Monte Carlo Policy Evaluation

Given a policy, $\pi$, estimate the value function, $V(s)$, for all states, $s \in \mathcal{S}$

Monte Carlo Policy Evaluation (first visit):

Input: a policy $\pi$ to be evaluated

Initialize:
    $V(s) \in \mathbb{R}$, arbitrarily, for all $s \in \mathcal{S}$
    $Returns(s) \leftarrow$ an empty list, for all $s \in \mathcal{S}$

Loop forever (for each episode):
    Generate an episode following $\pi$: $S_0, A_0, R_1, S_1, A_1, R_2, \ldots, S_{T-1}, A_{T-1}, R_T$
    $G \leftarrow 0$
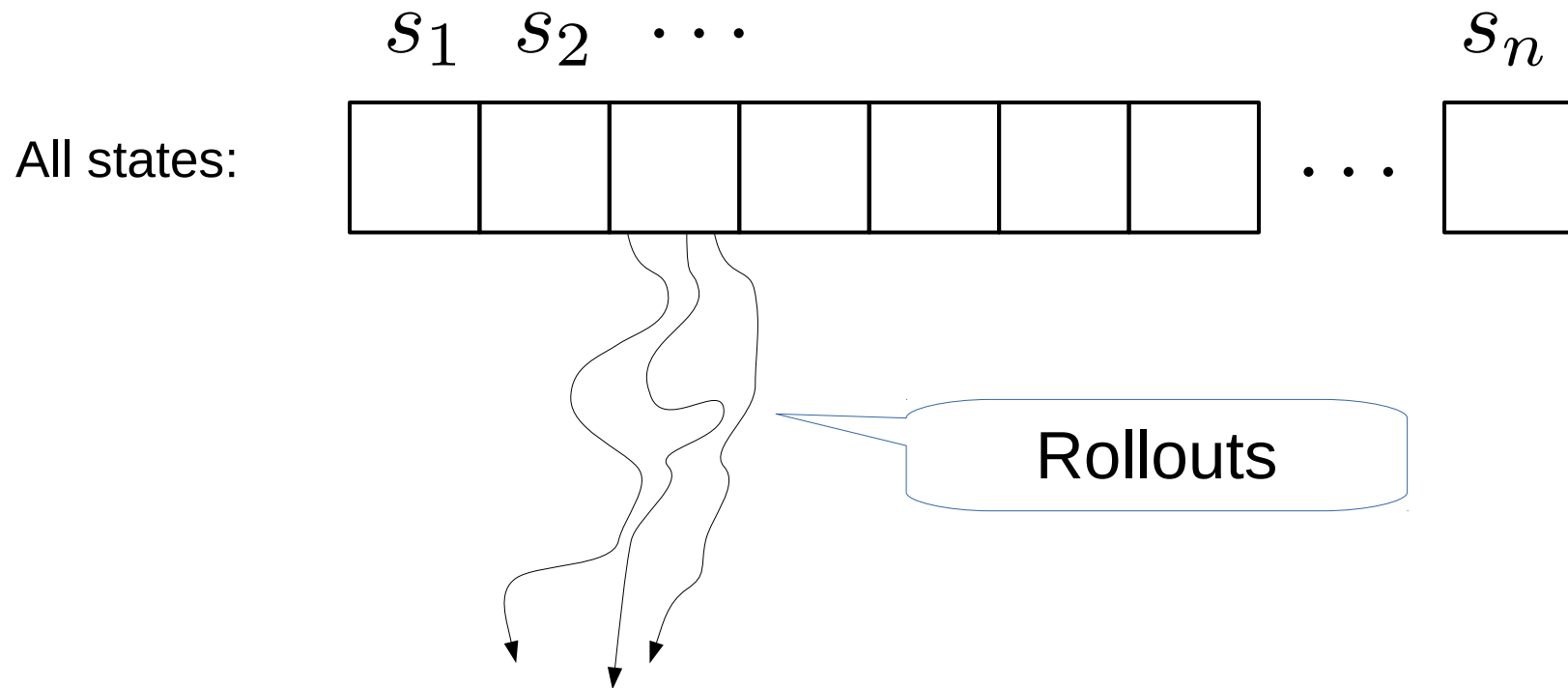    Loop for each step of episode, $t = T-1, T-2, \ldots, 0$:
        $G \leftarrow \gamma G + R_{t+1}$
        Unless $S_t$ appears in $S_0, S_1, \ldots, S_{t-1}$:
            Append $G$ to $Returns(S_t)$
            $V(S_t) \leftarrow average(Returns(S_t))$

# Monte Carlo Policy Evaluation

$$s_1 \quad s_2 \quad \cdots \qquad\qquad\qquad s_n$$

All states:



Rollouts

To get an accurate estimate of the value function, every state has to be visited many times.
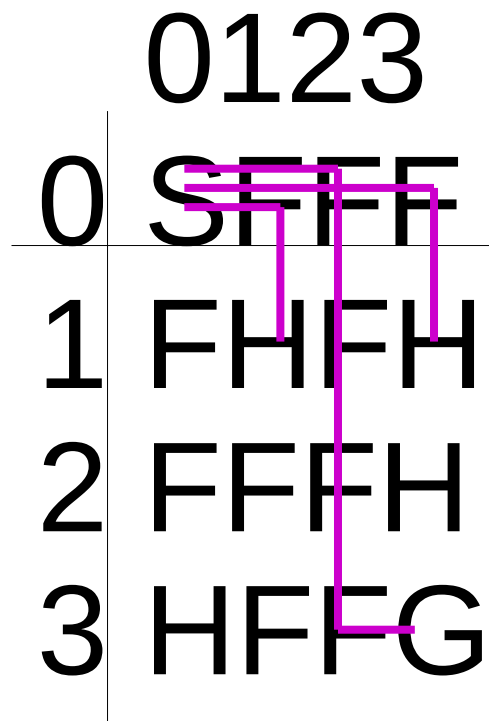
# Think-pair-share: frozenlake env

```
  0123
0 SFFF
1 FHFH
2 FFFH
3 HFFG
```

States: grid world coordinates

Actions: L, R, U, D

Reward: 0 except at G

# Think-pair-share: frozenlake env

0123

0 SFFF
1 FHFH
2 FFFH
3 HFFG

States: grid world coordinates

Actions: L, R, U, D
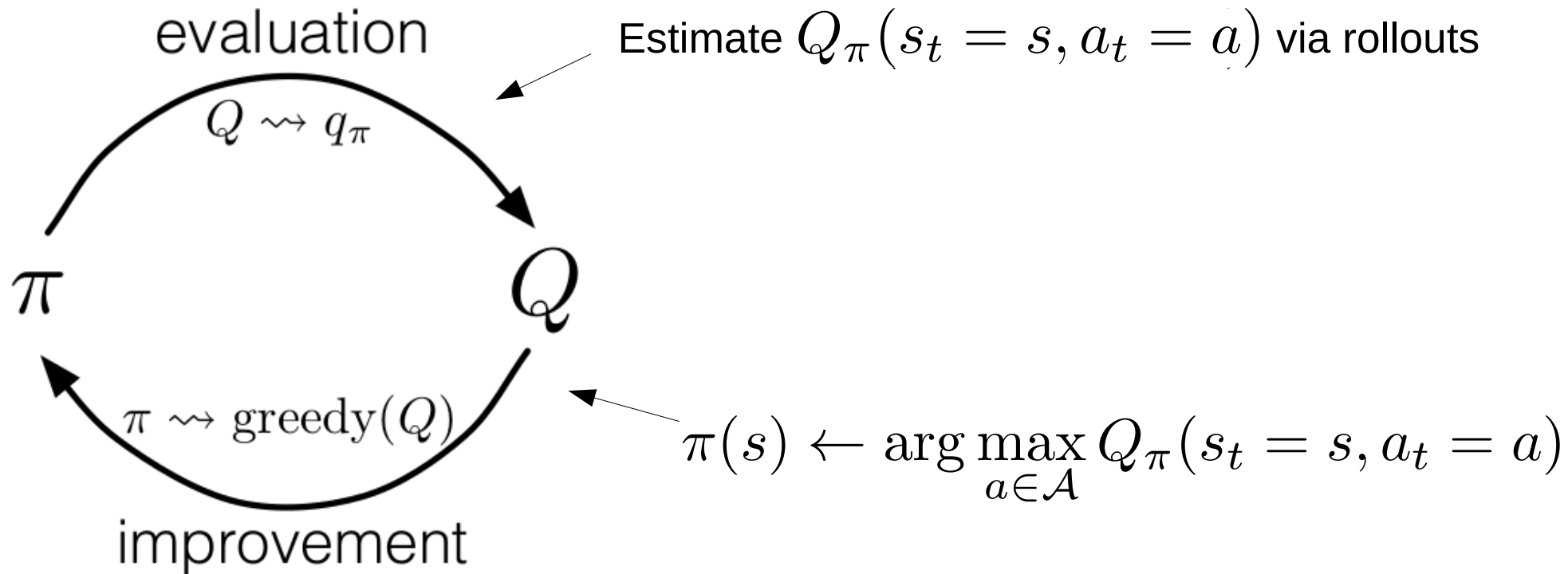
Reward: 0 except at G where r=1

Given: three episodes as shown

Calculate: values of states on top row as calculated by MC

# Monte Carlo Control

So far, we're only talking about policy *evaluation*

*…* but RL requires us to find a policy, not just evaluate it… How?



evaluation

$Q \rightsquigarrow q_\pi$

Estimate $Q_\pi(s_t = s, a_t = a)$ via rollouts

$\pi$

$Q$

$\pi \rightsquigarrow \mathrm{greedy}(Q)$

improvement

$\pi(s) \leftarrow \arg\max_{a \in \mathcal{A}} Q_\pi(s_t = s, a_t = a)$

Key idea: evaluate/improve policy iteratively...

# Monte Carlo Control

## Monte Carlo, Exploring Starts

Initialize:

$\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all $s \in \mathcal{S}$

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

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

Loop forever (for each episode):

Choose $S_0 \in \mathcal{S}$, $A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs have probability $> 0$

Generate an episode from $S_0, A_0$, following $\pi$: $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

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

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

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

Append $G$ to $Returns(S_t, A_t)$

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

$\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$

# Monte Carlo Control

Monte Carlo, Exploring Starts

Initialize:
$\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all $s$
$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s$
$Returns(s, a) \leftarrow$ empty list, for al.

Exploring starts:
– each episode starts with a random
   action taken from a random state

Loop forever (for each episode):
Choose $S_0 \in \mathcal{S}$, $A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs have probability $> 0$
Generate an episode from $S_0, A_0$, following $\pi$: $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$
$G \leftarrow 0$
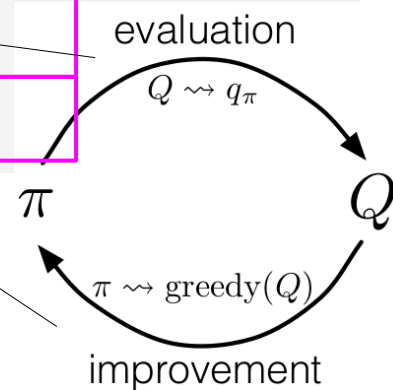Loop for each step of episode, $t = T-1, T-2, \ldots, 0$:
$G \leftarrow \gamma G + R_{t+1}$
Unless the pair $S_t, A_t$ appears in $S_0, A_0, S_1, A_1 \ldots, S_{t-1}, A_{t-1}$:
Append $G$ to $Returns(S_t, A_t)$
$Q(S_t, A_t) \leftarrow$ average$(Returns(S_t, A_t))$
$\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$

# Monte Carlo Control

## Monte Carlo, Exploring Starts

Initialize:
$\quad \pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all $s \in \mathcal{S}$
$\quad Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$
$\quad Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Loop forever (for each episode):
$\quad$ Choose $S_0 \in \mathcal{S}, A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs have probability $> 0$
$\quad$ Generate an episode from $S_0, A_0$, following $\pi$: $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$
$\quad G \leftarrow 0$
$\quad$ Loop for each step of episode, $t = T-1, T-2, \ldots, 0$:
$\quad\quad G \leftarrow \gamma G + R_{t+1}$
$\quad\quad$ Unless the pair $S_t, A_t$ appears in $S_0, A_0, S_1, A_1 \ldots, S_{t-1}, A_{t-1}$:
$\quad\quad\quad$ Append $G$ to $Returns(S_t, A_t)$
$\quad\quad\quad Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$
$\quad\quad\quad \pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$

evaluation

$Q \rightsquigarrow q_\pi$

$\pi$ $\qquad$ $Q$

$\pi \rightsquigarrow \text{greedy}(Q)$

improvement

# Monte Carlo Control

## Monte Carlo, Exploring Starts

Initialize:
  $\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all $s$
  $Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s$
  $Returns(s, a) \leftarrow$ empty list, for a

Loop forever (for each episode):
  Choose $S_0 \in \mathcal{S}$, $A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs ha$\quad$robability $> 0$
  Generate an episode from $S_0, A_0$, following $\pi$: $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$
  $G \leftarrow 0$
  Loop for each step of episode, $t = T-1, T-2, \ldots, 0$:
    $G \leftarrow \gamma G + R_{t+1}$
    Unless the pair $S_t, A_t$ appears in $S_0, A_0, S_1, A_1 \ldots, S_{t-1}, A_{t-1}$:
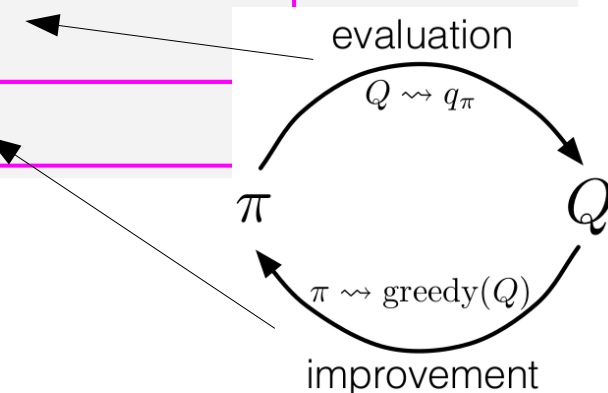      Append $G$ to $Returns(S_t, A_t)$
      $Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$
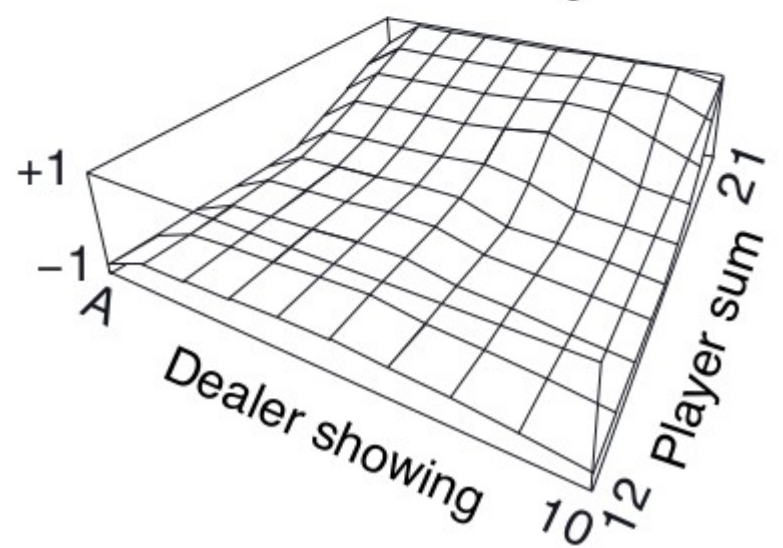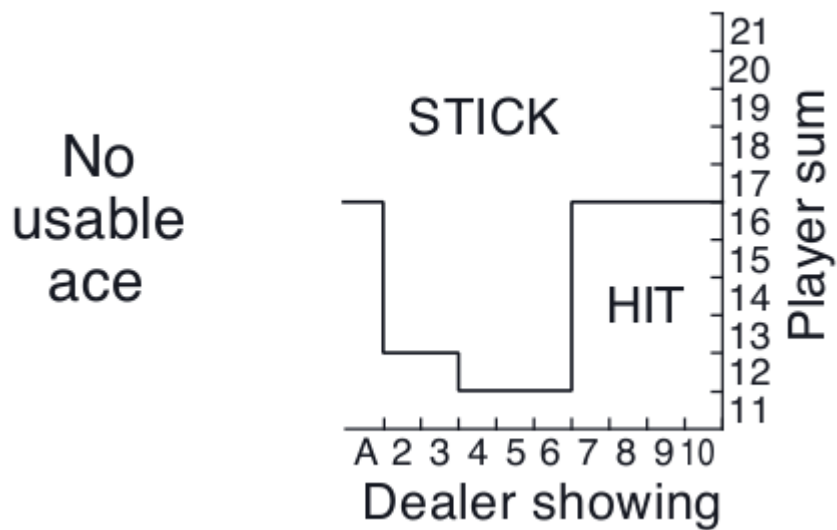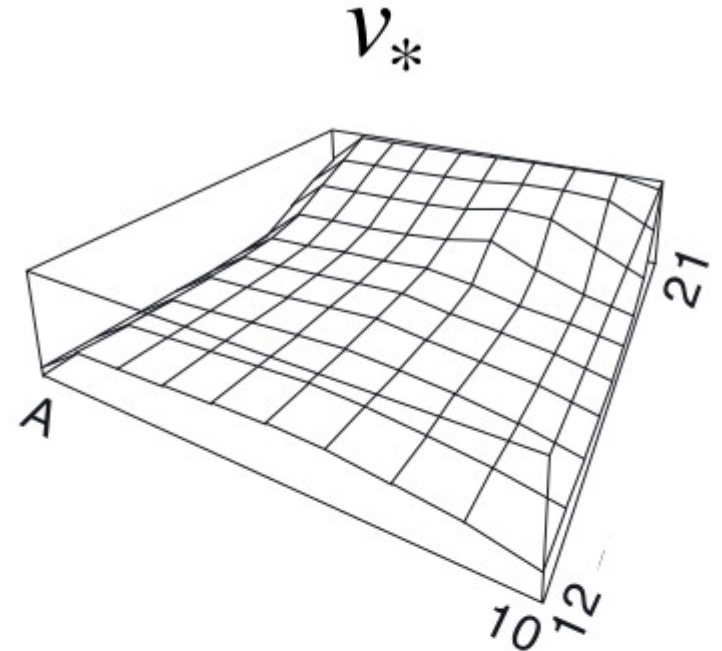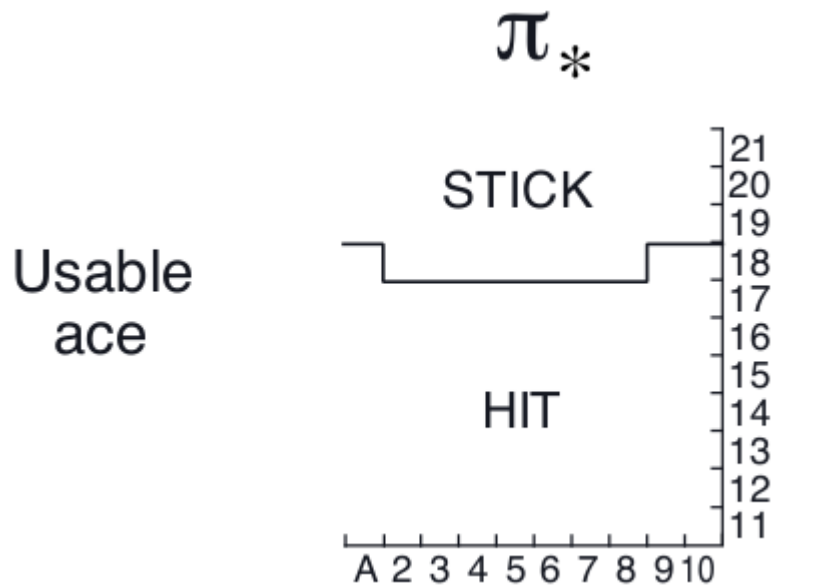      $\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$

Notice there is only one step of policy evaluation
 – that's okay.
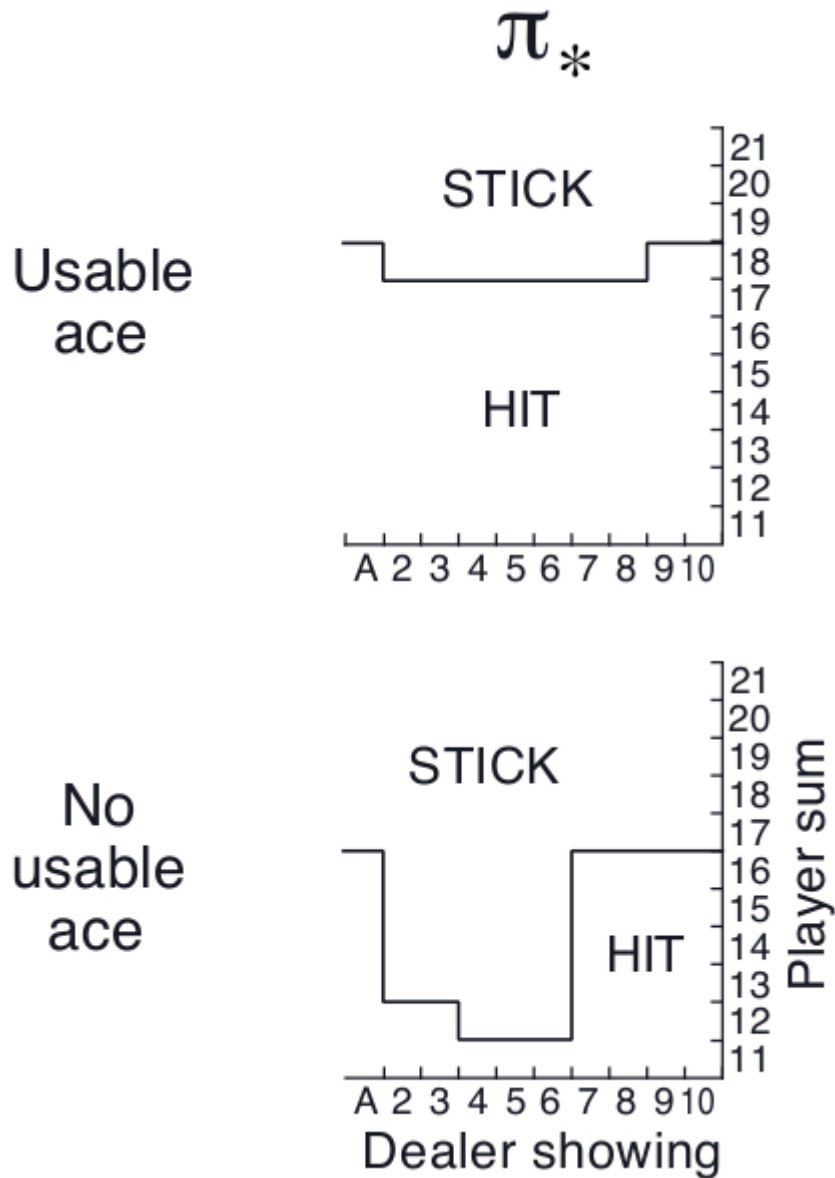 – each evaluation iter moves value fn toward its optimal value. Good enough to improve policy.
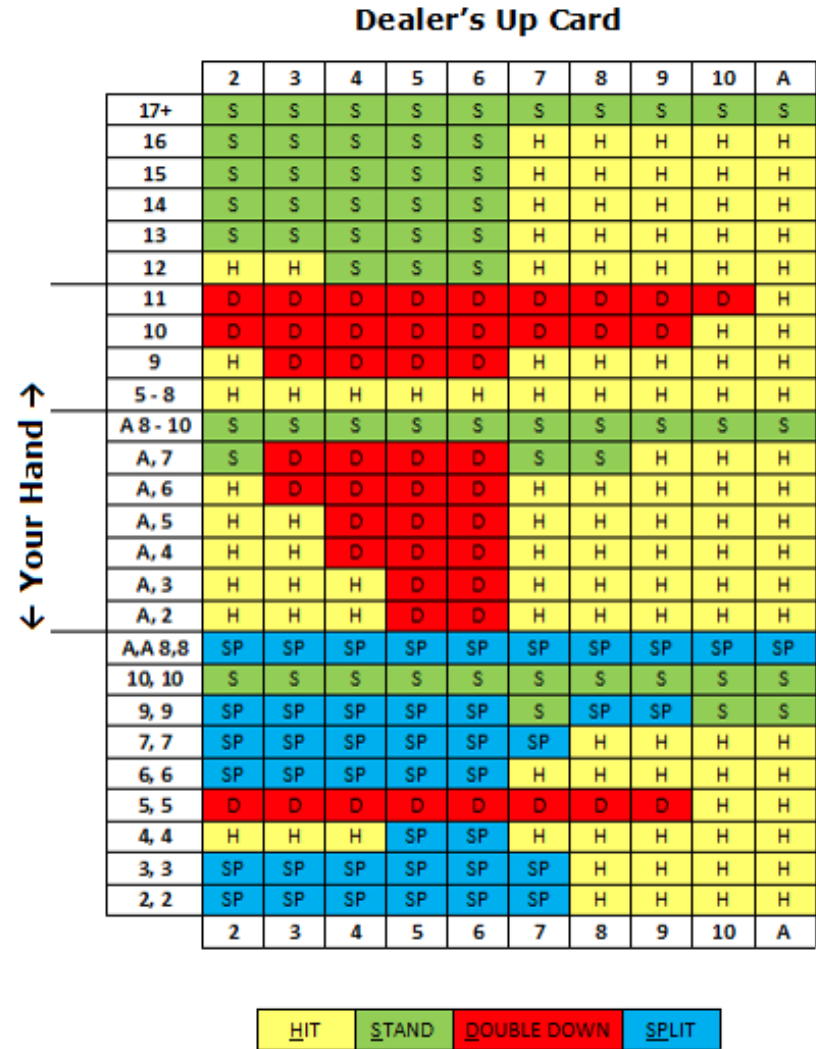
evaluation

$Q \rightsquigarrow q_\pi$

$\pi$ $\qquad$ $Q$

$\pi \rightsquigarrow \text{greedy}(Q)$

improvement

# Monte Carlo Control

# Monte Carlo Control



What the MC agent learned

The official "basic strategy"

☐ Greedified policy meets the conditions for policy improvement:

$$
\begin{aligned}
q_{\pi_k}(s, \pi_{k+1}(s)) &= q_{\pi_k}(s, \arg\max_a q_{\pi_k}(s, a)) \\
&= \max_a q_{\pi_k}(s, a) \\
&\geq q_{\pi_k}(s, \pi_k(s)) \\
&\geq v_{\pi_k}(s).
\end{aligned}
$$

☐ And thus must be $\geq \pi_k$ by the policy improvement theorem

☐ This assumes exploring starts and infinite number of episodes for MC policy evaluation

☐ To solve the latter:

- update only to a given level of performance
- alternate between evaluation and improvement per episode

# Monte Carlo Control: Convergence

❏ Greedified policy meets the conditions for policy improvement:

$$q_{\pi_k}(s, \pi_{k+1}(s)) = q_{\pi_k}(s, \arg\max_a q_{\pi_k}(s, a))$$

$$= \max_a q_{\pi_k}(s, a)$$

$$\geq q_{\pi_k}(s, \pi_k(s))$$

$$\geq v_{\pi_k}(s).$$

❏ And thus must be $\geq \pi_k$ by the policy improvement theorem

❏ This assumes exploring starts and in ̶ ̶ ̶ ̶ ̶ ̶ ̶ ̶ ̶ ̶ ̶ of episodes
for

❏ To

▪

If  $q_\pi(s, \pi'(s)) \geq v_\pi(s), \forall s \in S,$

then $v_{\pi'}(s) \geq v_\pi(s)$  i.e. $\pi'$ is better than $\pi$

▪ a ̶ ̶ ̶ ̶ ̶ ̶ ̶ ̶ ̶ isode

# Policy Improvement Theorem: Proof (Sketch)

$$v_\pi(s) \leq q_\pi(s, \pi'(s))$$

$$= \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = \pi'(s)] \qquad \text{(by (4.6))}$$

$$= \mathbb{E}_{\pi'}[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s]$$

$$\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma q_\pi(S_{t+1}, \pi'(S_{t+1})) \mid S_t = s] \qquad \text{(by (4.7))}$$

$$= \mathbb{E}_{\pi'}[R_{t+1} + \gamma \mathbb{E}_{\pi'}[R_{t+2} + \gamma v_\pi(S_{t+2}) \mid S_{t+1}, A_{t+1} = \pi'(S_{t+1})] \mid S_t = s]$$

$$= \mathbb{E}_{\pi'}\left[R_{t+1} + \gamma R_{t+2} + \gamma^2 v_\pi(S_{t+2}) \mid S_t = s\right]$$

$$\leq \mathbb{E}_{\pi'}\left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 v_\pi(S_{t+3}) \mid S_t = s\right]$$

$$\vdots$$

$$\leq \mathbb{E}_{\pi'}\left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \cdots \mid S_t = s\right]$$

$$= v_{\pi'}(s).$$

# E-Greedy Exploration

Monte Carlo, Exploring Starts:

Initialize:
$\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all $s$
$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s$
$Returns(s, a) \leftarrow$ empty list, for all

Loop forever (for each episode):
  Choose $S_0 \in \mathcal{S}$, $A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs have probability $> 0$
  Generate an episode from $S_0, A_0$, following $\pi$: $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$
  $G \leftarrow 0$
  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$(Returns(S_t, A_t))$
      $\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$

Without exploring starts, we are not guaranteed to explore the state/action space
– why is this a problem?
– what happens if we never experience certain transitions?

# E-Greedy Exploration

Monte Carlo, Exploring Starts:

Initialize:
$\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all $s$
$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s$
$Returns(s, a) \leftarrow$ empty list, for all $s$

Loop forever (for each episode):
Choose $S_0 \in \mathcal{S}$, $A_0 \in \mathcal{A}(S_0)$ randomly such
Generate an episode from $S_0, A_0$, following
$G \leftarrow 0$
Loop for each step of episode, $t = T-1, T-$
$G \leftarrow \gamma G + R_{t+1}$
Unless the pair $S_t, A_t$ appears in $S_0, A_0, S_1, A_1 \ldots, S_{t-1}, A_{t-1}$:
Append $G$ to $Returns(S_t, A_t)$
$Q(S_t, A_t) \leftarrow$ average($Returns(S_t, A_t)$)
$\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$

Without exploring starts, we are not guaranteed to explore the state/action space
– why is this a problem?
– what happens if we never experience certain transitions?

Can we accomplish this without exploring starts?

# E-Greedy Exploration

Monte Carlo, Exploring Starts:

Initialize:
$\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all $s$
$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s$
$Returns(s, a) \leftarrow$ empty list, for all

Loop forever (for each episode):
Choose $S_0 \in \mathcal{S}$, $A_0 \in \mathcal{A}(S_0)$ randomly such
Generate an episode from $S_0, A_0$, followin
$G \leftarrow 0$
Loop for each step of episode, $t = T-1, T$
$G \leftarrow \gamma G + R_{t+1}$
Unless the pair $S_t, A_t$ appears in $S_0$
Append $G$ to $Returns(S_t, A_t)$
$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$
$\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$

Without exploring starts, we are not guaranteed to explore the state/action space
– why is this a problem?
– what happens if we never experience certain transitions?

Can we accomplish this without exploring starts?

Yes: create a stochastic (e-greedy) policy

# E-Greedy Exploration

Greedy policy:

$$\pi(a|s) = \begin{cases} 1 & \text{if } a = a^* \\ 0 & \text{otherwise} \end{cases} \qquad a^* = \arg\max_{a \in \mathcal{A}} Q(s, a)$$

E-Greedy policy:

$$\pi(a|s) = \begin{cases} 1 - \epsilon + \dfrac{\epsilon}{|\mathcal{A}|} & \text{if } a = a^* \\[2ex] \dfrac{\epsilon}{|\mathcal{A}|} & \text{otherwise} \end{cases}$$

# E-Greedy Exploration

Greedy policy:

$$\pi(a|s) = \begin{cases} 1 & \text{if } a = a^* \\ 0 & \text{otherwise} \end{cases} \qquad a^* = \arg\max_{a \in \mathcal{A}} Q(s,a)$$

E-Greedy policy:

$$\pi(a|s) = \begin{cases} 1 - \epsilon + \dfrac{\epsilon}{|\mathcal{A}|} & \text{if } a = a^* \\ \dfrac{\epsilon}{|\mathcal{A}|} & \text{otherwise} \end{cases}$$

Action drawn uniformly from $\mathcal{A}$

# E-Greedy Exploration

Greedy policy:

$$\pi(a|s) = \begin{cases} 1 \\ \\ 0 \end{cases}$$

Guarantees every state/action will be visited infinitely often

E-Greedy policy:

$$\pi(a|s) = \begin{cases} 1 - \epsilon + \dfrac{\epsilon}{|\mathcal{A}|} & \text{if } a = a^* \\ \\ \dfrac{\epsilon}{|\mathcal{A}|} & \text{otherwise} \end{cases}$$

– Notice that this is a stochastic policy (not deterministic).
– This is an example of an *soft* policy
– *soft policy*: all actions in all states have non-zero probability

# E-Greedy Exploration

Monte Carlo, ε-greedy exploration:

Algorithm parameter: small $\varepsilon > 0$

Initialize:

    $\pi \leftarrow$ an arbitrary $\varepsilon$-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 $\pi$: $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$

    $G \leftarrow 0$

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

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

        Unless the pair $S_t, A_t$ appears in $S_0, A_0, S_1, A_1$ ...

            Append $G$ to $Returns(S_t, A_t)$

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

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

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

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

E-greedy exploration

# Off-Policy Methods

- *On-policy* methods evaluate or improve the policy that is used to make decisions.
- *Off-policy* methods evaluate or improve a policy different from that used to generate the data.

- The *target policy* is the policy (π) we wish to evaluate/improve.
- The *behavior policy* is the policy (b) used to generate experiences.

- Coverage:

$$\forall s, a \left[ \pi(a|s) > 0 \implies b(a|s) > 0 \right]$$

# MC Summary

MC methods estimate value function by doing rollouts

Can estimate either the state value function, $V(s)$, or the action value function, $Q(s, a)$

MC Control alternates between policy evaluation and policy improvement

E-greedy exploration explores all possible actions while preferring greedy actions

Off-policy methods update a policy other than the one used to generate experience