

# Chapter 5: Monte Carlo Methods

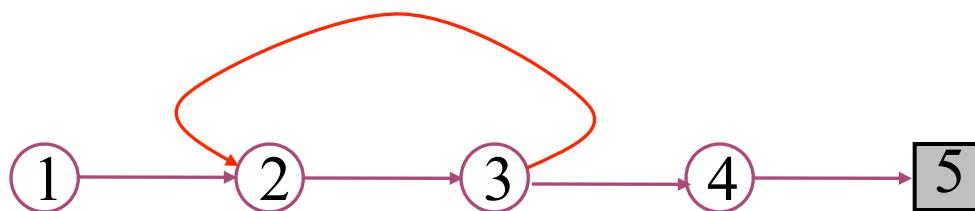
---

- Monte Carlo methods are learning methods  
Experience → values, policy
- Monte Carlo methods can be used in two ways:
  - *model-free*: No model necessary and still attains optimality
  - *Simulated*: Needs only a simulation, not a *full* model
- Monte Carlo methods learn from *complete* sample returns
  - Only defined for episodic tasks (in this book)
- Like an associative version of a bandit method

# Monte Carlo Policy Evaluation

---

- ❑ *Goal*: learn  $v_\pi(s)$
- ❑ *Given*: some number of episodes under  $\pi$  which contain  $s$
- ❑ *Idea*: Average returns observed after visits to  $s$



- ❑ *Every-Visit MC*: average returns for *every* time  $s$  is visited in an episode
- ❑ *First-visit MC*: average returns only for *first* time  $s$  is visited in an episode
- ❑ Both converge asymptotically

# First-visit Monte Carlo policy evaluation

---

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:

    Generate an episode using  $\pi$

    For each state  $s$  appearing in the episode:

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

        Append  $G$  to  $Returns(s)$

$V(s) \leftarrow$  average( $Returns(s)$ )

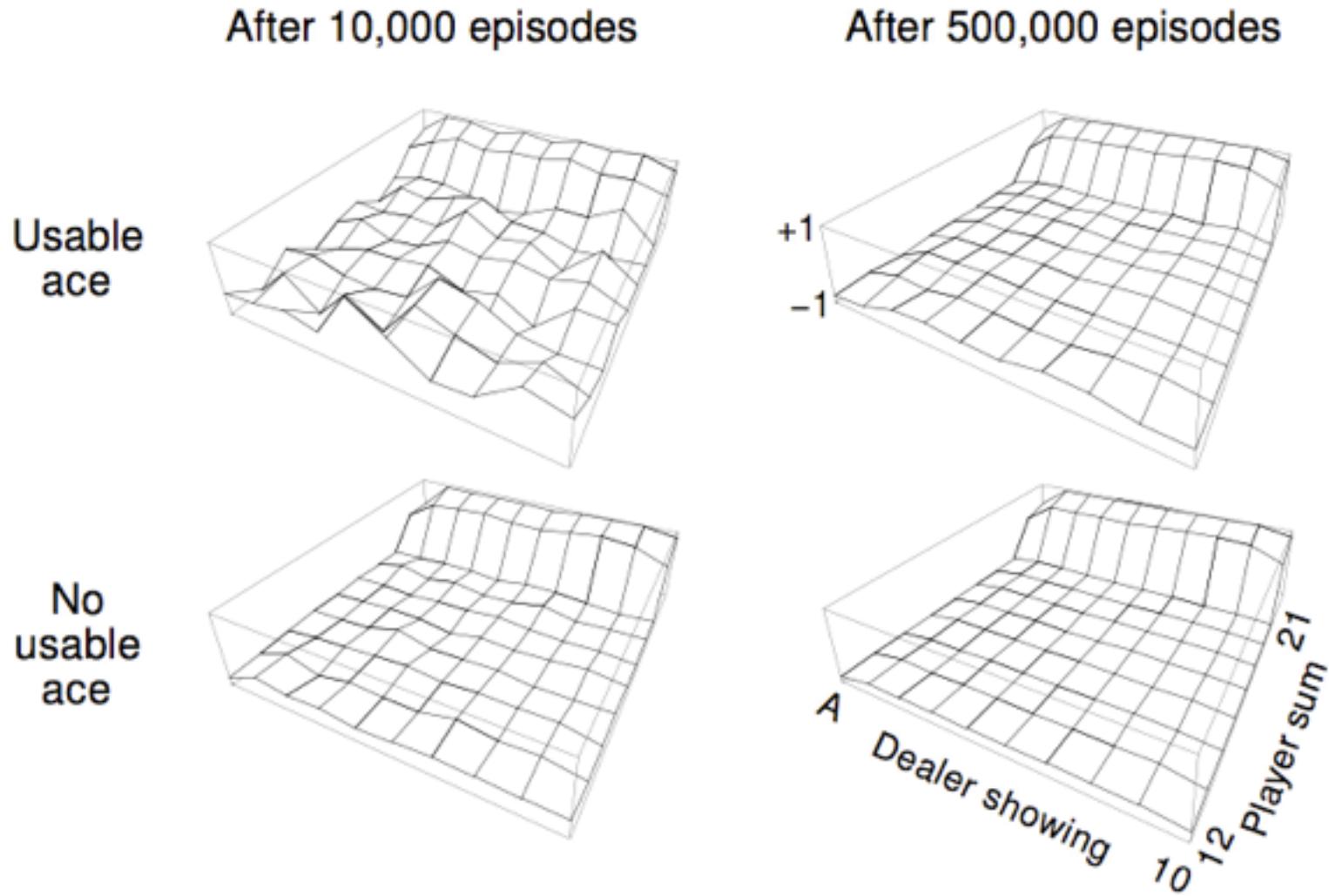
# Blackjack example

---

- *Object*: Have your card sum be greater than the dealer's without exceeding 21.
- *States* (200 of them):
  - current sum (12-21)
  - dealer's showing card (ace-10)
  - do I have a useable ace?
- *Reward*: +1 for winning, 0 for a draw, -1 for losing
- *Actions*: stick (stop receiving cards), hit (receive another card)
- *Policy*: Stick if my sum is 20 or 21, else hit
- No discounting ( $\gamma = 1$ )



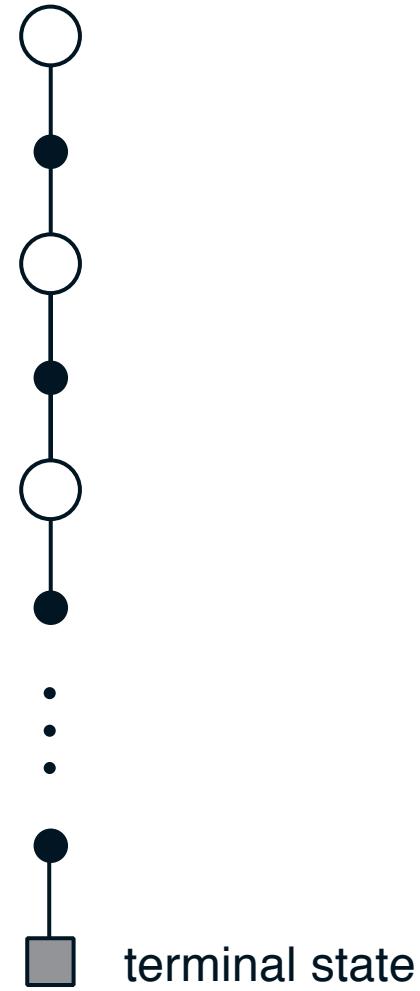
# Learned blackjack state-value functions



# Backup diagram for Monte Carlo

---

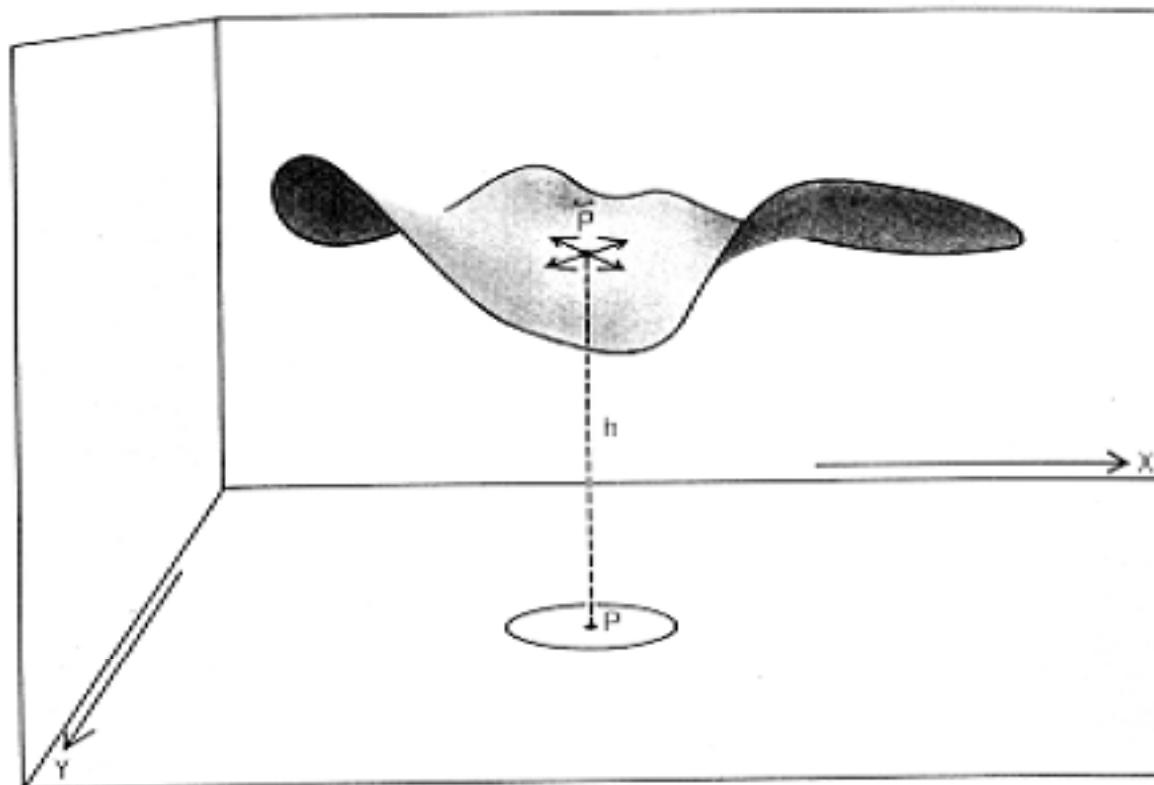
- ❑ Entire rest of episode included
- ❑ Only one choice considered at each state (unlike DP)
  - thus, there will be an explore/exploit dilemma
- ❑ Does not bootstrap from successor states's values (unlike DP)
- ❑ Time required to estimate one state does not depend on the total number of states



# The Power of Monte Carlo

e.g., Elastic Membrane (Dirichlet Problem)

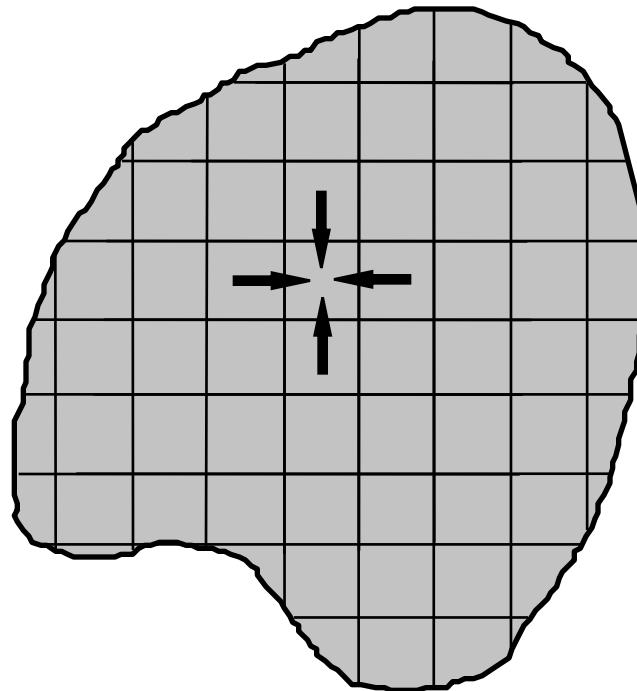
How do we compute the shape of the membrane or bubble?



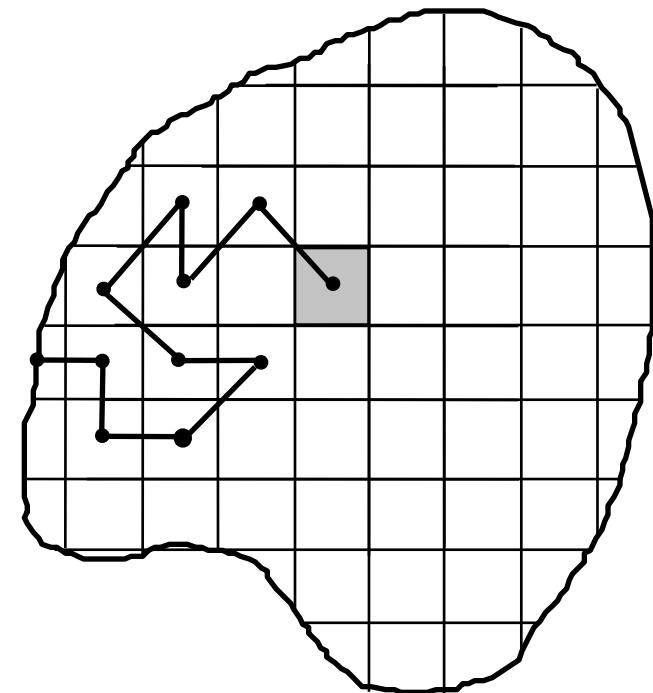
# Two Approaches

---

Relaxation



Kakutani's algorithm, 1945



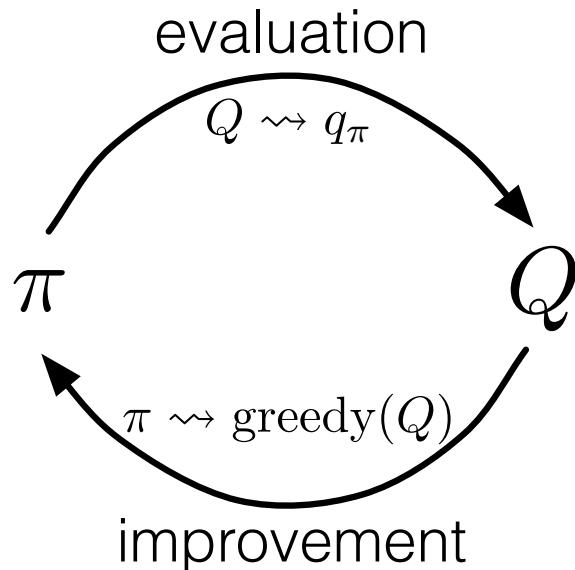
# Monte Carlo Estimation of Action Values (Q)

---

- Monte Carlo is most useful when a model is not available
  - We want to learn  $q^*$
- $q_\pi(s,a)$  - average return starting from state  $s$  and action  $a$  following  $\pi$
- Converges asymptotically *if* every state-action pair is visited
- *Exploring starts*: Every state-action pair has a non-zero probability of being the starting pair

# Monte Carlo Control

---



- MC policy iteration: Policy evaluation using MC methods followed by policy improvement
- Policy improvement step: greedify with respect to value (or action-value) function

# Convergence of MC Control

---

- 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 Exploring Starts

---

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

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

$$\pi(s) \leftarrow \text{arbitrary}$$

$$Returns(s, a) \leftarrow \text{empty list}$$

Fixed point is optimal policy  $\pi^*$

Now proven (almost)

Repeat forever:

Choose  $S_0 \in \mathcal{S}$  and  $A_0 \in \mathcal{A}(S_0)$  s.t. all pairs have probability  $> 0$

Generate an episode starting from  $S_0, A_0$ , following  $\pi$

For each pair  $s, a$  appearing in the episode:

$G \leftarrow$  return following the first occurrence of  $s, a$

Append  $G$  to  $Returns(s, a)$

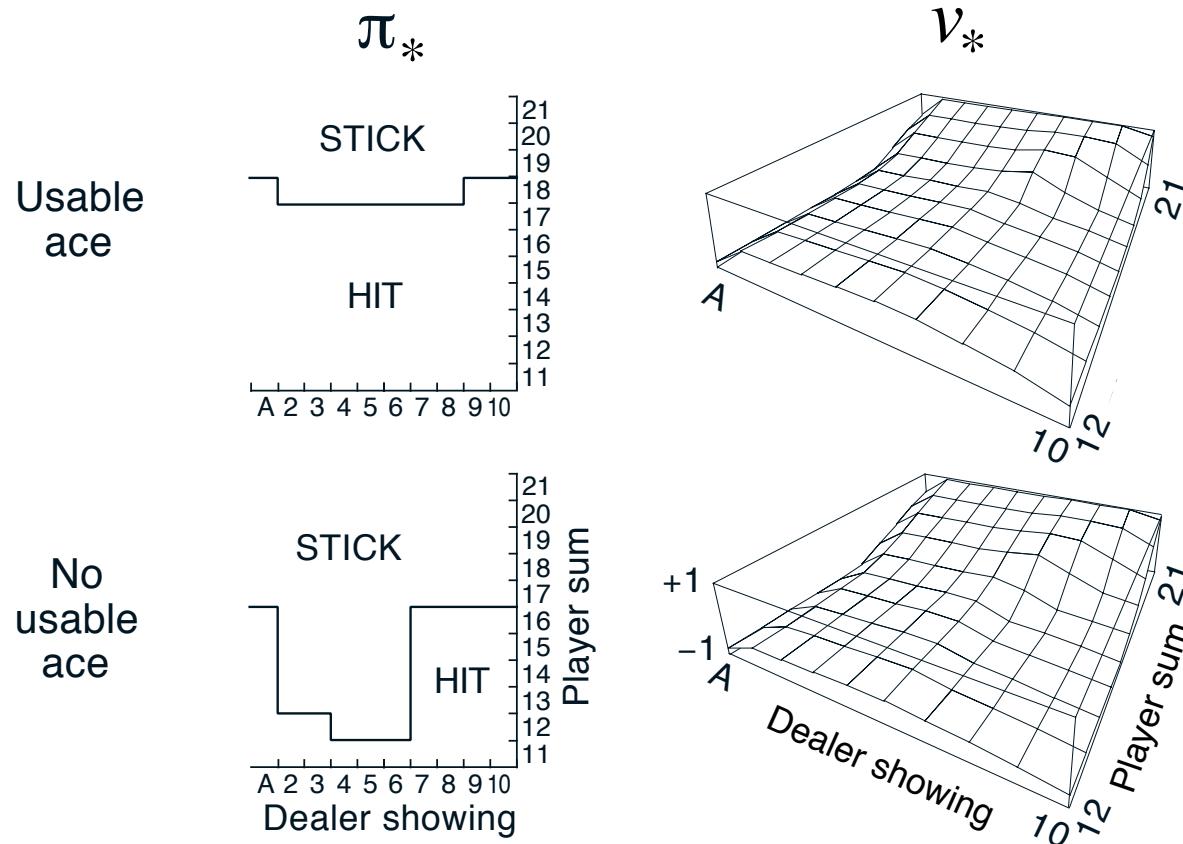
$Q(s, a) \leftarrow \text{average}(Returns(s, a))$

For each  $s$  in the episode:

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

# Blackjack example continued

- ❑ Exploring starts
- ❑ Initial policy as described before



# On-policy Monte Carlo Control

---

- *On-policy*: learn about policy currently executing
- How do we get rid of exploring starts?
  - The policy must be eternally *soft*:
    - $\pi(a|s) > 0$  for all  $s$  and  $a$
  - e.g.  $\epsilon$ -soft policy:
    - probability of an action = 
$$\begin{array}{ll} \frac{\epsilon}{|\mathcal{A}(s)|} & \text{or} \\ \text{non-max} & 1 - \epsilon + \frac{\epsilon}{|\mathcal{A}(s)|} \\ & \max \text{ (greedy)} \end{array}$$
- Similar to GPI: move policy *towards* greedy policy (e.g.,  $\epsilon$ -greedy)
- Converges to best  $\epsilon$ -soft policy

# On-policy MC Control

---

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

$Q(s, a) \leftarrow$  arbitrary

$Returns(s, a) \leftarrow$  empty list

$\pi(a|s) \leftarrow$  an arbitrary  $\varepsilon$ -soft policy

Repeat forever:

(a) Generate an episode using  $\pi$

(b) For each pair  $s, a$  appearing in the episode:

$G \leftarrow$  return following the first occurrence of  $s, a$

Append  $G$  to  $Returns(s, a)$

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

(c) For each  $s$  in the episode:

$A^* \leftarrow \arg \max_a Q(s, a)$

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

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

# What we've learned about Monte Carlo so far

---

- MC has several advantages over DP:
  - Can learn directly from interaction with environment
  - No need for full models
  - No need to learn about ALL states (no bootstrapping)
  - Less harmed by violating Markov property (later in book)
- MC methods provide an alternate policy evaluation process
- One issue to watch for: maintaining sufficient exploration
  - exploring starts, soft policies

# Off-policy methods

---

- Learn the value of the *target policy*  $\pi$  from experience due to *behavior policy*  $\mu$
- For example,  $\pi$  is the greedy policy (and ultimately the optimal policy) while  $\mu$  is exploratory (e.g.,  $\varepsilon$ -soft)
- In general, we only require *coverage*, i.e., that  $\mu$  generates behavior that covers, or includes,  $\pi$

$$\pi(a|s) > 0 \quad \text{for every } s, a \text{ at which} \quad \mu(a|s) > 0$$

- Idea: *importance sampling*
  - Weight each return by the *ratio of the probabilities* of the trajectory under the two policies

# Importance Sampling Ratio

---

- Probability of the rest of the trajectory, after  $S_t$ , under  $\pi$ :

$$\begin{aligned} & \Pr\{A_t, S_{t+1}, A_{t+1}, \dots, S_T \mid S_t, A_{t:T-1} \sim \pi\} \\ &= \pi(A_t|S_t)p(S_{t+1}|S_t, A_t)\pi(A_{t+1}|S_{t+1}) \cdots p(S_T|S_{T-1}, A_{T-1}) \\ &= \prod_{k=t}^{T-1} \pi(A_k|S_k)p(S_{k+1}|S_k, A_k), \end{aligned}$$

- In importance sampling, each return is weighted by the relative probability of the trajectory under the two policies

$$\rho_t^T = \frac{\prod_{k=t}^{T-1} \pi(A_k|S_k)p(S_{k+1}|S_k, A_k)}{\prod_{k=t}^{T-1} \mu(A_k|S_k)p(S_{k+1}|S_k, A_k)} = \prod_{k=t}^{T-1} \frac{\pi(A_k|S_k)}{\mu(A_k|S_k)}$$

- This is called the *importance sampling ratio*
- All importance sampling ratios have expected value 1

$$\mathbb{E}_{A_k \sim \mu} \left[ \frac{\pi(A_k|S_k)}{\mu(A_k|S_k)} \right] = \sum_a \mu(a|S_k) \frac{\pi(a|S_k)}{\mu(a|S_k)} = \sum_a \pi(a|S_k) = 1.$$

# Importance Sampling

---

- New notation: time steps increase across episode boundaries:

- $\dots \textcolor{red}{s} \dots \textcolor{green}{\#} \dots \dots \textcolor{black}{\#} \dots \dots \textcolor{red}{s} \dots \dots \textcolor{green}{\#} \dots \dots$
- $t = 1 \ 2 \ 3 \ \textcolor{red}{4} \ 5 \ 6 \ 7 \ 8 \ \textcolor{green}{9} \ 10 \ 11 \ 12 \ 13 \ 14 \ 15 \ 16 \ 17 \ 18 \ 19 \ \textcolor{red}{20} \ 21 \ 22 \ 23 \ 24 \ \textcolor{green}{25} \ 26 \ 27$



$\mathcal{T}(s) = \{4, 20\}$   
set of start times

$T(4) = 9 \quad T(20) = 25$   
next termination times

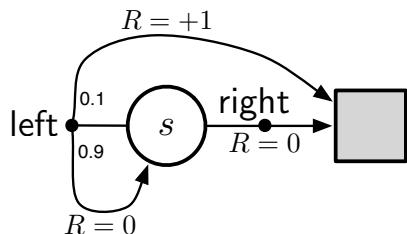
- *Ordinary importance sampling* forms estimate

$$V(s) \doteq \frac{\sum_{t \in \mathcal{T}(s)} \rho_t^{T(t)} G_t}{|\mathcal{T}(s)|}$$

- Whereas *weighted importance sampling* forms estimate

$$V(s) \doteq \frac{\sum_{t \in \mathcal{T}(s)} \rho_t^{T(t)} G_t}{\sum_{t \in \mathcal{T}(s)} \rho_t^{T(t)}}$$

# Example of infinite variance under *ordinary* importance sampling



$$\pi(\text{left}|s) = 1$$

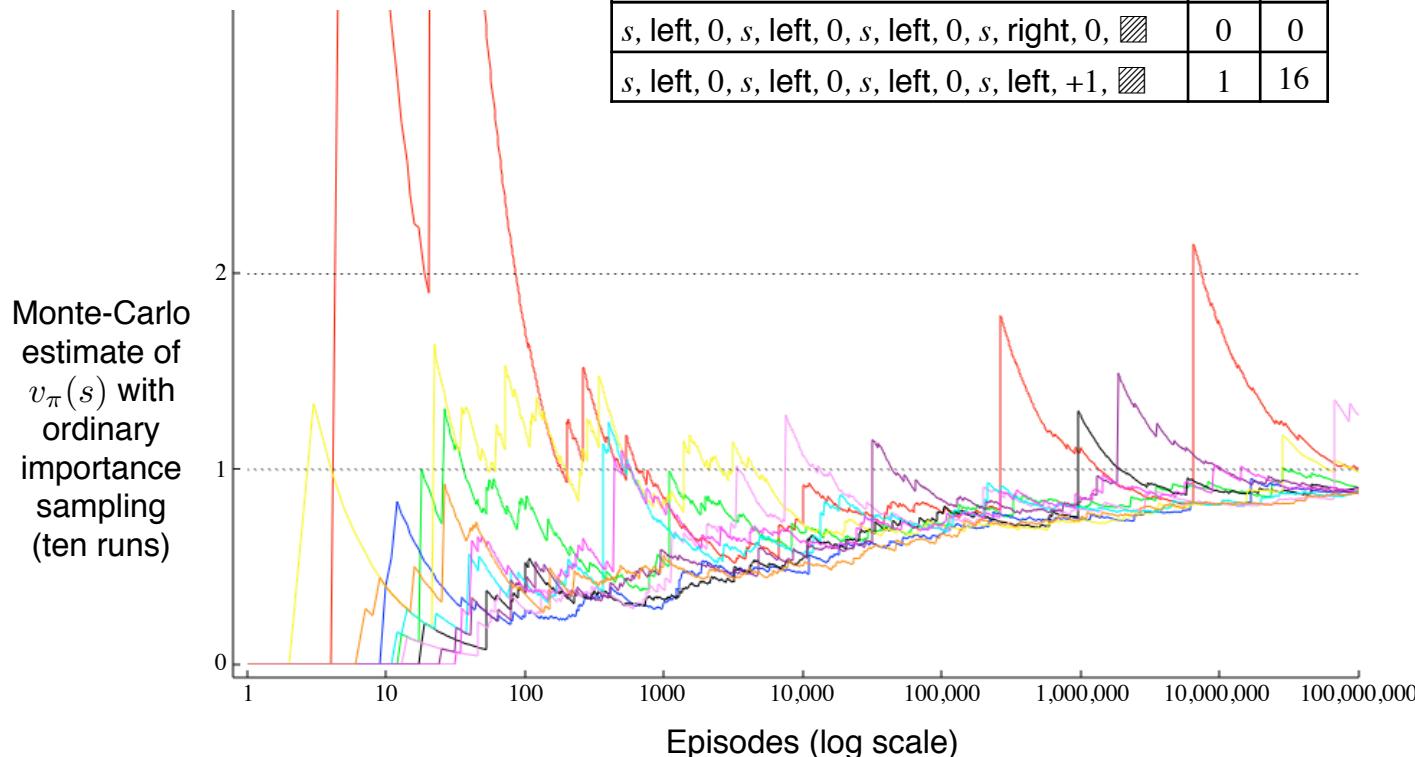
$$\mu(\text{left}|s) = \frac{1}{2}$$

$$\gamma = 1$$

$$v_\pi(s) = 1$$

$$\frac{\pi(\text{right}|s)}{\mu(\text{right}|s)} = 0$$

$$\frac{\pi(\text{left}|s)}{\mu(\text{left}|s)} = 2$$



OIS:

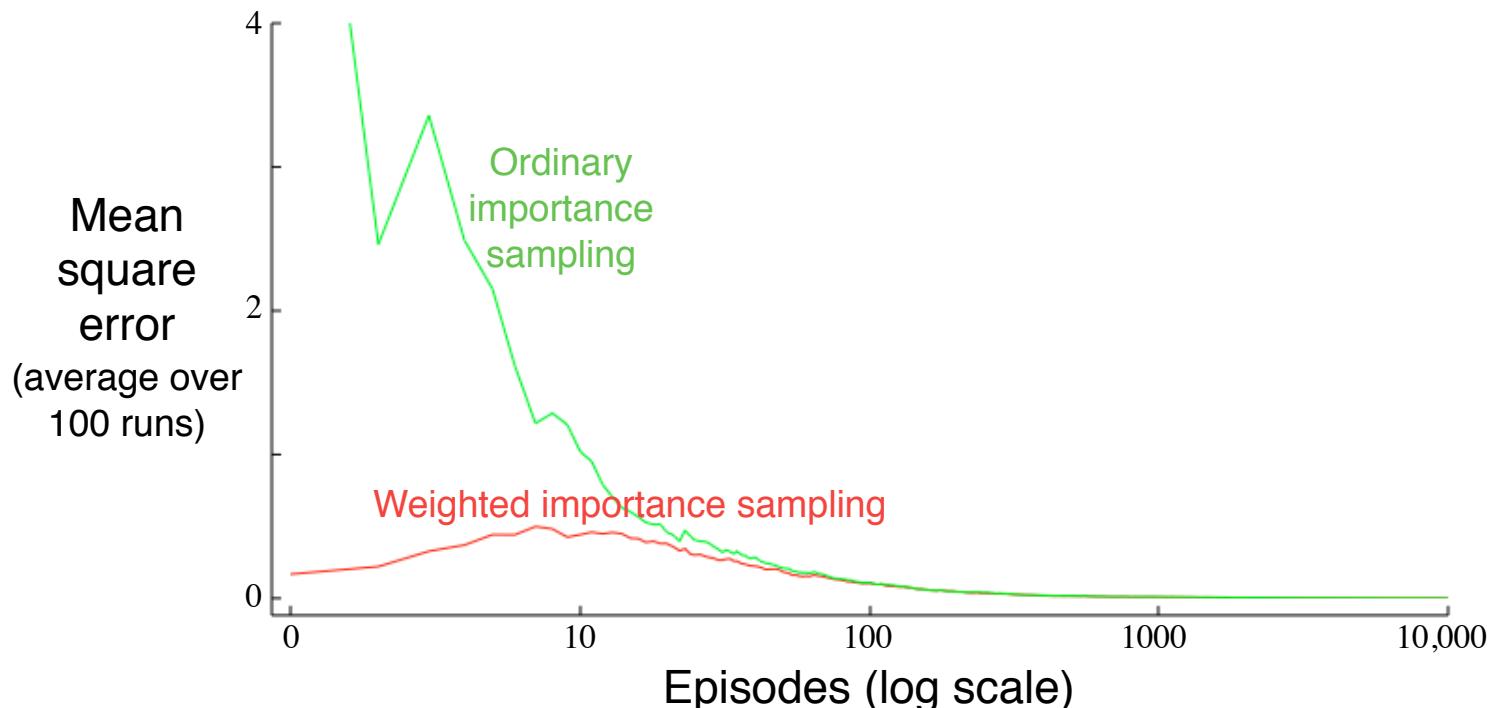
$$V(s) \triangleq \frac{\sum_{t \in \mathcal{T}(s)} \rho_t^{T(t)} G_t}{|\mathcal{T}(s)|}$$

WIS:

$$V(s) \triangleq \frac{\sum_{t \in \mathcal{T}(s)} \rho_t^{T(t)} G_t}{\sum_{t \in \mathcal{T}(s)} \rho_t^{T(t)}}$$

# Example: Off-policy Estimation of the value of a *single* Blackjack State

- State is player-sum 13, dealer-showing 2, useable ace
- Target policy is stick only on 20 or 21
- Behavior policy is equiprobable
- True value  $\approx -0.27726$



## Incremental off-policy every-visit MC policy evaluation (returns $Q \approx q_\pi$ )

Input: an arbitrary target policy  $\pi$

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

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

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

Repeat forever:

$\mu \leftarrow$  any policy with coverage of  $\pi$

Generate an episode using  $\mu$ :

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

$$G \leftarrow 0$$

$$W \leftarrow 1$$

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

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

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

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

$$W \leftarrow W \frac{\pi(A_t | S_t)}{\mu(A_t | S_t)}$$

If  $W = 0$  then ExitForLoop

## Off-policy every-visit MC control (returns $\pi \approx \pi_*$ )

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:

$\mu \leftarrow$  any soft policy

Generate an episode using  $\mu$ :

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

$$G \leftarrow 0$$

$$W \leftarrow 1$$

For  $t = T - 1, T - 2, \dots$  downto 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 ExitForLoop

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

Target policy is greedy  
and deterministic

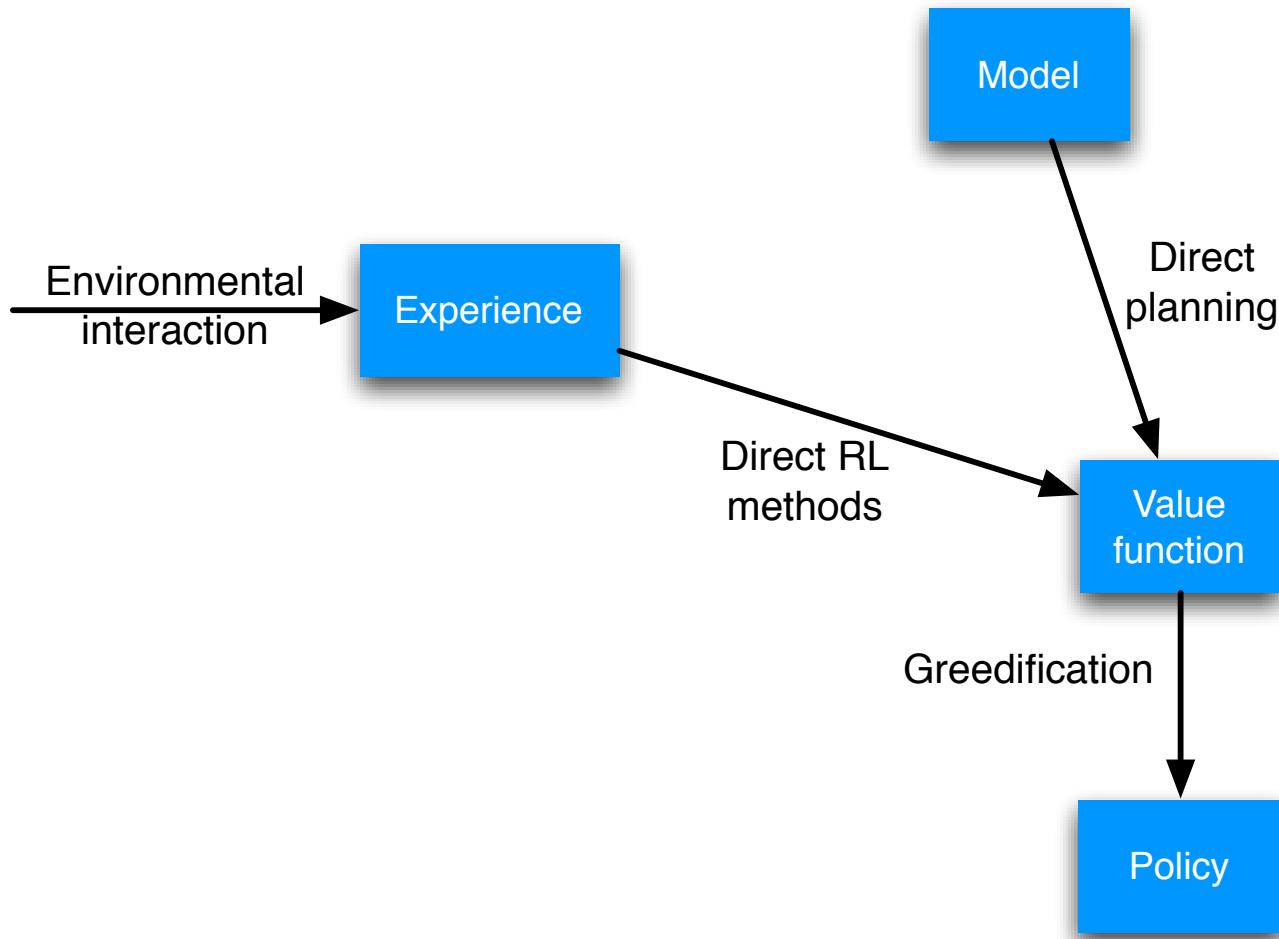
Behavior policy is soft,  
typically  $\varepsilon$ -greedy

# Summary

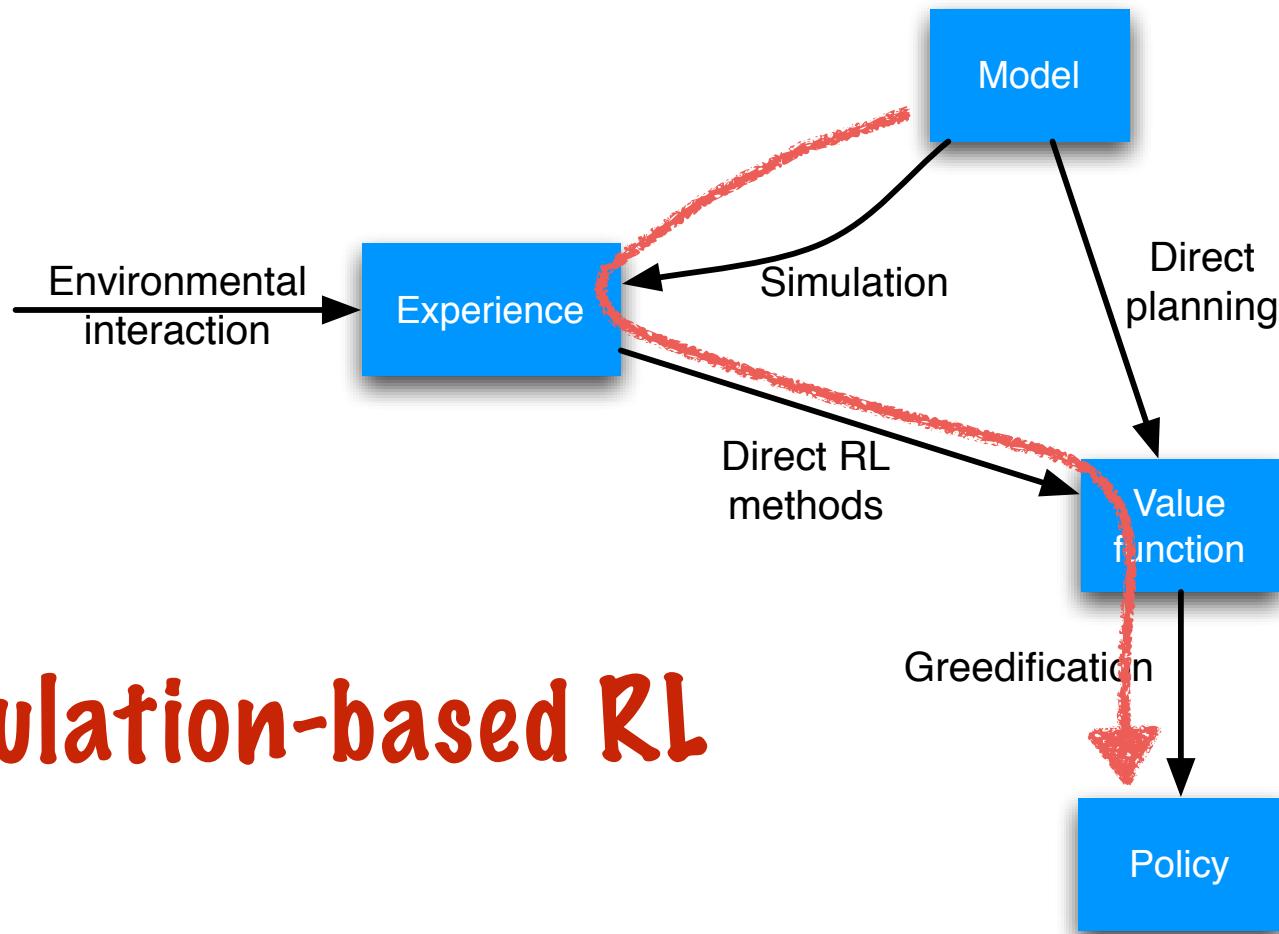
---

- ❑ MC has several advantages over DP:
  - Can learn directly from interaction with environment
  - No need for full models
  - Less harmed by violating Markov property (later in book)
- ❑ MC methods provide an alternate policy evaluation process
- ❑ One issue to watch for: maintaining sufficient exploration
  - exploring starts, soft policies
- ❑ Introduced distinction between *on-policy* and *off-policy* methods
- ❑ Introduced *importance sampling* for off-policy learning
- ❑ Introduced distinction between *ordinary* and *weighted* IS

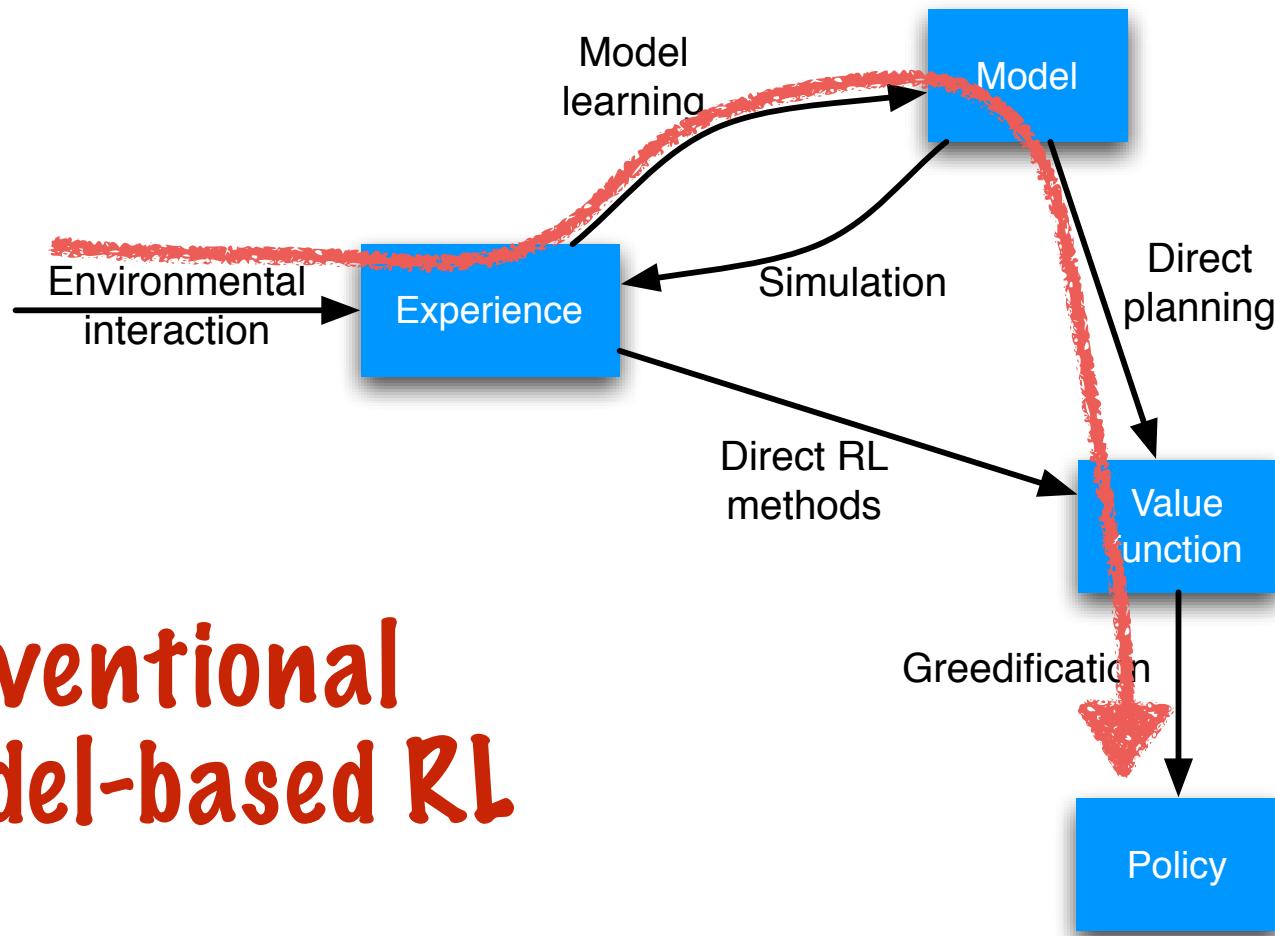
# Paths to a policy



# Paths to a policy

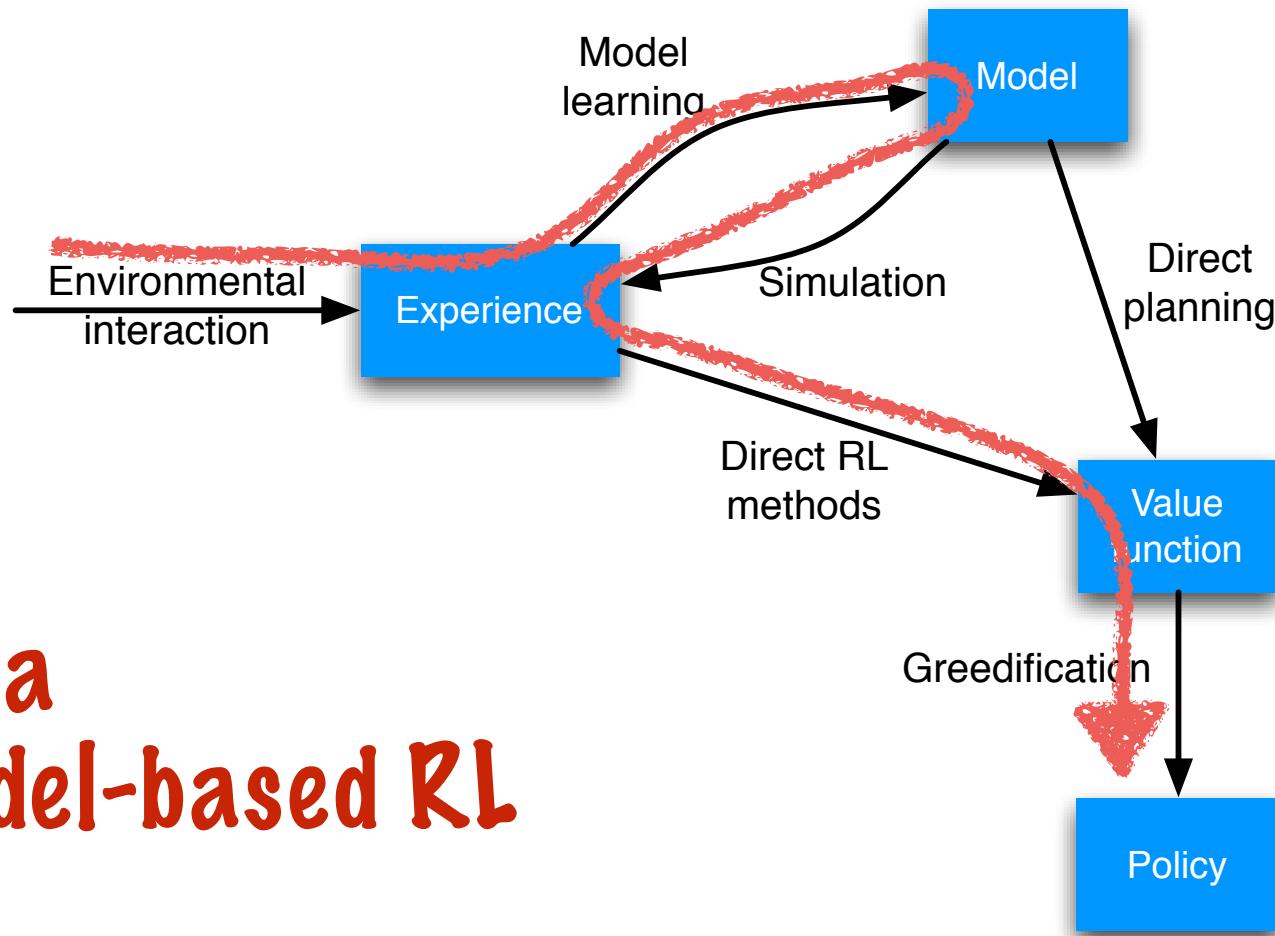


# Paths to a policy



**Conventional  
Model-based RL**

# Paths to a policy



Dyna  
Model-based RL