



Reinforcement Learning

Computer Engineering Department

Sharif University of Technology

Mohammad Hossein Rohban, Ph.D.

Spring 2025

Courtesy: Some slides are adopted from CS 285 Berkeley, and CS 234 Stanford, and Pieter Abbeel's compact series on RL.

Value Function

$$\mathbb{E}(g(x)) = \sum_{x_i \in \mathcal{D}(x)} P(x=x_i) g(x_i)$$

- $V^*(s)$ = expected utility **starting in s** , and **acting optimally** in all subsequent actions.

Bellman's Opt.
Eq.

$$V^*(s) = \max_{a_i^s} \mathbb{E} \left(\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t, s_{t+1}) \mid s_0 = s \right)$$

$\underbrace{R(s_0, a_0, s_1) + \gamma \sum_{t=1}^{\infty} \gamma^{t-1} R(s_t, a_t, s_{t+1})}_{\gamma V^*(s')}$

$$= \max_{a_0} \sum_{s_1=s'} P(s'|s_0, a_0) [R(s_0, a_0, s') + \gamma V^*(s')]$$

Value Iteration

- $\underline{V_0^*}(s)$ = optimal value for state s when $H=0$
 - $V_0^*(s) = 0 \quad \forall s$
- $V_1^*(s)$ = optimal value for state s when $H=1$
 - $V_1^*(s) = \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_0^*(s'))$
- $V_2^*(s)$ = optimal value for state s when $H=2$
 - $V_2^*(s) = \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_1^*(s'))$
- $V_k^*(s)$ = optimal value for state s when $H = k$
 - $V_k^*(s) = \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}^*(s'))$

Value Iteration

Algorithm:

Start with $V_0^*(s) = 0$ for all s .

For $k = 1, \dots, H$:

For all states s in S :

$$V_k^*(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}^*(s'))$$

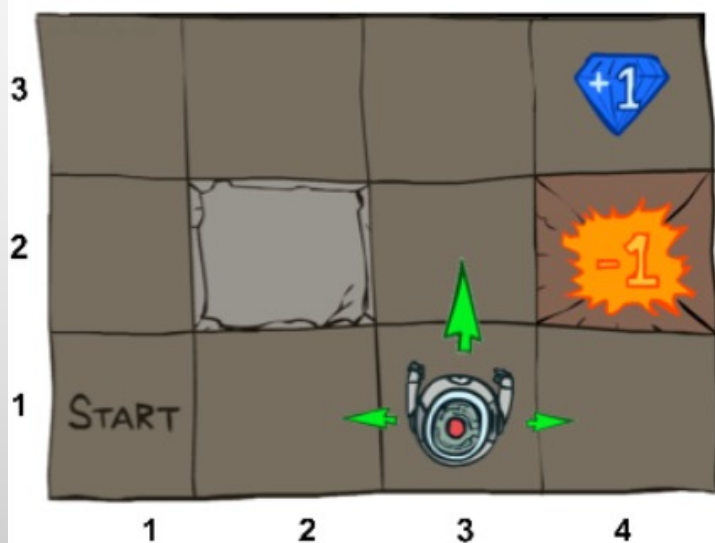
$$\pi_k^*(s) \leftarrow \arg \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}^*(s'))$$

This is called a **value update** or **Bellman update/back-up**

Value Iteration

$$\underline{V_0(s)} \leftarrow 0$$

$k = 0$



Noise = 0.2

Discount = 0.9

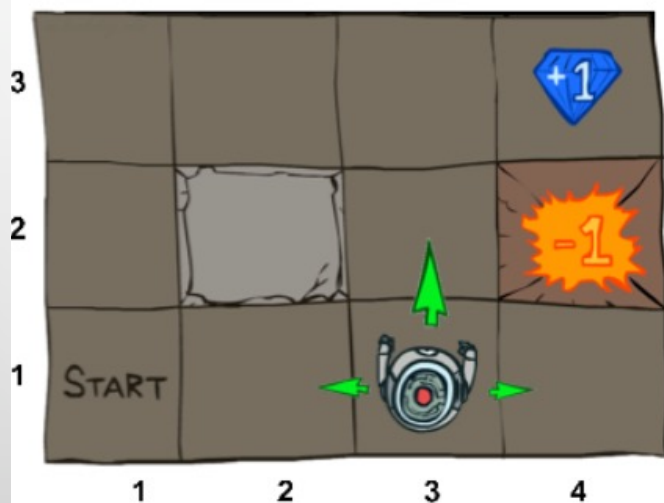
0.00	0.00	0.00	0.00
0.00		0.00	0.00
0.00	0.00	0.00	0.00

VALUES AFTER 0 ITERATIONS

Value Iteration

$$V_1(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_0(s'))$$

k = 0



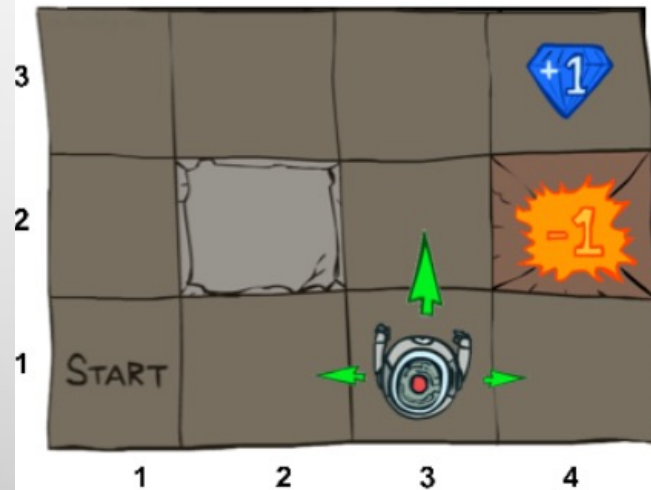
Noise = 0.2
Discount = 0.9

0.00	0.00	0.00	0.00
0.00		0.00	0.00
0.00	0.00	0.00	0.00
VALUES AFTER 0 ITERATIONS			

Value Iteration

$$V_2(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_1(s'))$$

k = 1



Noise = 0.2
Discount = 0.9

0.00	0.00	0.00	1.00
0.00		0.00	-1.00
0.00	0.00	0.00	0.00

VALUES AFTER 1 ITERATIONS

Value Iteration

$$V_2(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_1(s'))$$

k = 2

0.00	0.00	0.72	1.00
0.00		0.00	-1.00
0.00	0.00	0.00	0.00

VALUES AFTER 2 ITERATIONS

Noise = 0.2
Discount = 0.9

Value Iteration

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_k(s'))$$

k = 3

0.00	0.52	0.78	1.00
0.00		0.43	-1.00
0.00	0.00	0.00	0.00

VALUES AFTER 3 ITERATIONS

Noise = 0.2

Discount = 0.9

Value Iteration

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_k(s'))$$

k = 4

0.37	0.66	0.83	1.00
0.00		0.51	-1.00
0.00	0.00	0.31	0.00

VALUES AFTER 4 ITERATIONS

Noise = 0.2

Discount = 0.9

Value Iteration

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_k(s'))$$

k = 5

0.51	0.72	0.84	1.00
0.27		0.55	-1.00
0.00	0.22	0.37	0.13

VALUES AFTER 5 ITERATIONS

Noise = 0.2
Discount = 0.9

Value Iteration

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_k(s'))$$

k = 6

0.59	0.73	0.85	1.00
0.41		0.57	-1.00
0.21	0.31	0.43	0.19

VALUES AFTER 6 ITERATIONS

Noise = 0.2
Discount = 0.9

Value Iteration

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_k(s'))$$

k = 7

0.62	0.74	0.85	1.00
0.50		0.57	-1.00
0.34	0.36	0.45	0.24

VALUES AFTER 7 ITERATIONS

Noise = 0.2
Discount = 0.9

Value Iteration

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_k(s'))$$

k = 8

0.63	0.74	0.85	1.00
0.53		0.57	-1.00
0.42	0.39	0.46	0.26

VALUES AFTER 8 ITERATIONS

Noise = 0.2
Discount = 0.9

Value Iteration

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_k(s'))$$

k = 9



Noise = 0.2
Discount = 0.9

Value Iteration

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_k(s'))$$

k = 10

0.64	0.74	0.85	1.00
0.56		0.57	-1.00
0.48	0.41	0.47	0.27

VALUES AFTER 10 ITERATIONS

Noise = 0.2
Discount = 0.9

Value Iteration

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_k(s'))$$

k = 11

0.64	0.74	0.85	1.00
0.56		0.57	-1.00
0.48	0.42	0.47	0.27

VALUES AFTER 11 ITERATIONS

Noise = 0.2
Discount = 0.9

Value Iteration

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_k(s'))$$

k = 12

0.64	0.74	0.85	1.00
0.57		0.57	-1.00
0.49	0.42	0.47	0.28

VALUES AFTER 12 ITERATIONS

Noise = 0.2
Discount = 0.9

Value Iteration

$$V_{k+1}(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_k(s'))$$

k = 100

0.64	0.74	0.85	1.00
0.57		0.57	-1.00
0.49	0.43	0.48	0.28

VALUES AFTER 100 ITERATIONS

Noise = 0.2

Discount = 0.9

Q-Values

- $Q^*(s, a)$ = expected utility **starting in s , taking action a , and (thereafter) acting optimally**

$$V^*(s) = \max_{a'} Q^*(s, a')$$

$$\max_{a'} \underbrace{Q^*(s', a')}$$

$$Q^*(s, a) = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V^*(s')]$$

- Bellman Equation:

$$Q^*(s, a) = \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma \max_{a'} Q^*(s', a'))$$

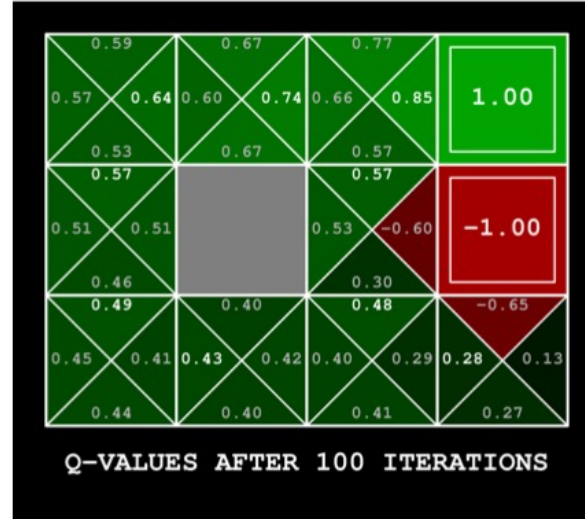
- Q-value Iteration:

$$Q_{k+1}^*(s, a) \leftarrow \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma \max_{a'} Q_k^*(s', a'))$$

Q-Value Iteration

$$Q_{k+1}^*(s, a) \leftarrow \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma \max_{a'} Q_k^*(s', a'))$$

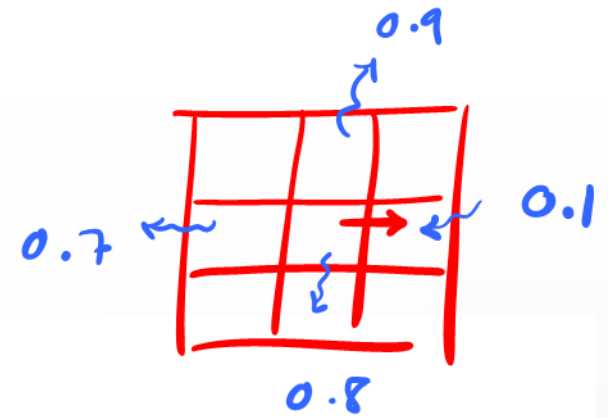
k = 100



Noise = 0.2

Discount = 0.9

Policy Evaluation



- Recall value iteration:

$$V_k^*(s) \leftarrow \max_a \sum_{s'} P(s'|s, a) (R(s, a, s') + \gamma V_{k-1}^*(s'))$$

- Policy evaluation for a given $\pi(s)$:

$$V_k^\pi(s) \leftarrow \sum_{s'} P(s'|s, \pi(s)) (R(s, \pi(s), s') + \gamma V_{k-1}^\pi(s'))$$

At convergence:

$$\forall s \quad V^\pi(s) \leftarrow \sum_{s'} P(s'|s, \pi(s)) (R(s, \pi(s), s') + \gamma V^\pi(s))$$

Policy Iteration

- One iteration of policy iteration

- Policy evaluation for current policy π_k :

- Iterate until convergence

Eval

$$\rightarrow \underline{V_{i+1}^{\pi_k}(s)} \leftarrow \sum_{s'} P(s'|s, \pi_k(s)) [R(s, \pi(s), s') + \gamma \underline{V_i^{\pi_k}(s')}]$$

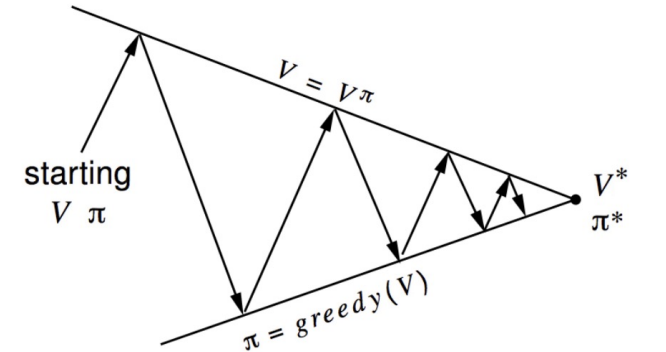
- Policy improvement: find the best action according to one-step look-ahead

Improvement

$$\rightarrow \pi_{k+1}(s) \leftarrow \arg \max_a \sum_{s'} \underbrace{P(s'|s, a)} \underbrace{[R(s, a, s') + \gamma \underline{V_i^{\pi_k}(s')}]}$$

- Repeat until policy converges

- At convergence: optimal policy; and converges faster than value iteration under some conditions



One-step look ahead improves the policy

- Consider an alternative policy $\pi_{(k+1)}^{(1)}(t, s)$ that takes the prescribed actions in $\pi_{k+1}(s)$ **only at time $t = 0$** ; and stays the same as $\pi_k(s)$ in later times.
- The value function $V(s)$ for this new **time-dependent** policy is larger than or equal to $V(s)$ for the original policy $\pi_k(s)$ for all s . Why?
- Now let $\pi_{(k+1)}^{(2)}(t, s)$, which takes the prescribed action in $\pi_{k+1}(s)$ **only at times $t = 0$ and $t = 1$** , and stays the same as $\pi_k(s)$ in later times.
- Similarly, $V(s)$ gets improved for $\pi_{(k+1)}^{(2)}(t, s)$ compared to $\pi_{(k+1)}^{(1)}(t, s)$ for all s .
- Repeating this argument $\pi_{(k+1)}^{(\infty)}(t, s)$ becomes the same as $\pi_{k+1}(s)$.

An Example

Let this be the initial policy π_0 , show how policy improvement, makes this policy better.

+1	→	+1
-1	↑	-1
-1	→	-1
-1	↑	-1
-1	←	-1
-1	↑	-1
-1	→	-1

Planning vs. Learning

$$E(g(x)) \approx \frac{1}{n} \sum_{i=1}^n g(x_i)$$

$$x_1 \dots x_n \stackrel{iid}{\sim} X$$

- Assumed to have access to the dynamics $P(s' | s, a)$.
- We **don't have access** to this in the real world.
- We need to **estimate (or learn)** the value functions.

$$Q^*(s, a) = \sum_{s'} P(s' | s, a) (R(s, a, s') + \gamma \max_{a'} Q^*(s', a'))$$

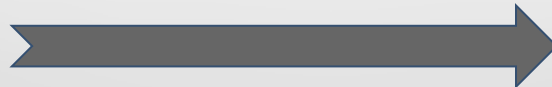
Monte-Carlo Prediction

Monte Carlo Methods - Introduction

- **Experience** samples to learn without a model
- MC methods require only experience— sample sequences of states, actions, and rewards from actual or simulated interaction with an environment.
- We can learn with samples: **episodes**!

We don't have access to

$$P(s'|s, a)$$



Model Free Learning!

Monte-Carlo prediction

- Suppose we wish to estimate $V_\pi(s)$, the value of a state s under policy π .
- The **first-visit mc** method estimates $V_\pi(s)$ as *the average of the returns* following first visits to s .

First-visit MC prediction, for estimating $V \approx v_\pi$

Input: a policy π 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 π : $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

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

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

Unless S_t appears in S_0, S_1, \dots, S_{t-1} :

Append G to $Returns(S_t)$

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

Episodes: another example

Input Policy π

	A	
B	C	D
	E	

Assume: $\gamma = 1$

Observed Episodes (Training)

Episode 1

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

Episode 2

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

Episode 3

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

Episode 4

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

Output Values

	-10	
+8	+4	+10
	-2	

Every Visit Monte-Carlo Policy

Initialize $N(s) = 0$, $G(s) = 0 \forall s \in S$

Loop

- Sample episode $i = s_{i,1}, a_{i,1}, r_{i,1}, s_{i,2}, a_{i,2}, r_{i,2}, \dots, s_{i,T_i}$
- Define $G_{i,t} = r_{i,t} + \gamma r_{i,t+1} + \gamma^2 r_{i,t+2} + \dots + \gamma^{T_i-t} r_{i,T_i}$ as return from time step t onwards in i th episode
- For each time step t until T_i (the end of the episode i)
 - state s is the state visited at time step t in episodes i
 - Increment counter of total visits: $N(s) = N(s) + 1$
 - Increment total return $G(s) = G(s) + G_{i,t}$
 - Update estimate $V^\pi(s) = G(s)/N(s)$

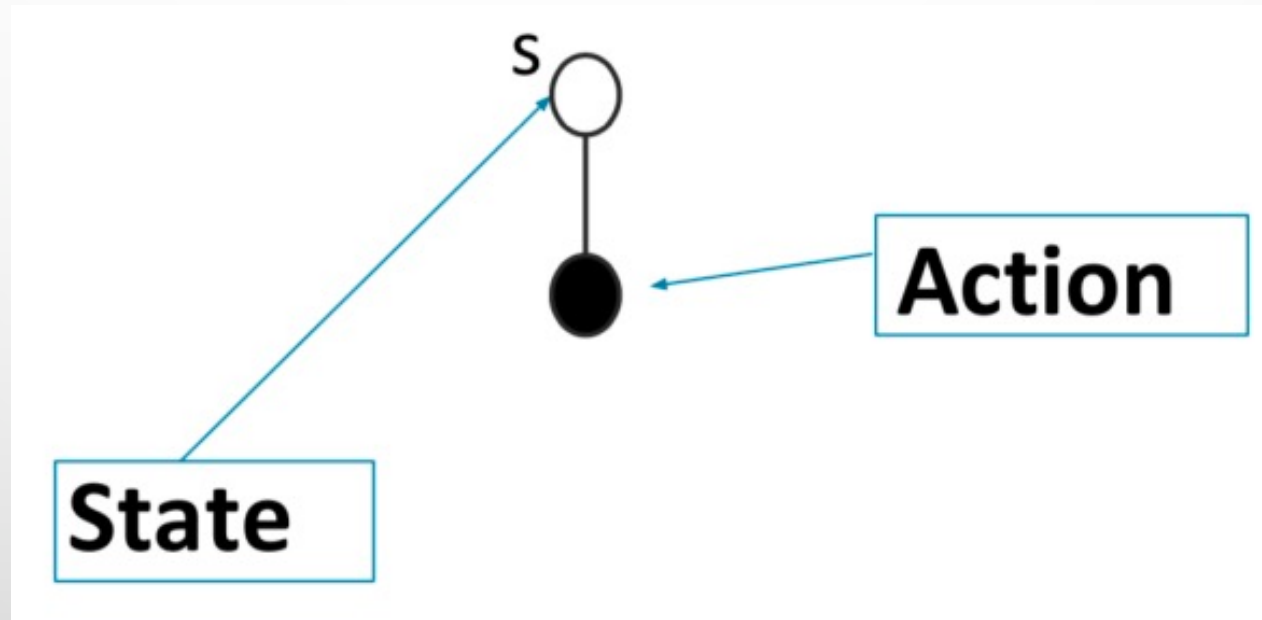
Incremental Monte-Carlo Policy

After each episode $i = s_{i,1}, a_{i,1}, r_{i,1}, s_{i,2}, a_{i,2}, r_{i,2}, \dots$

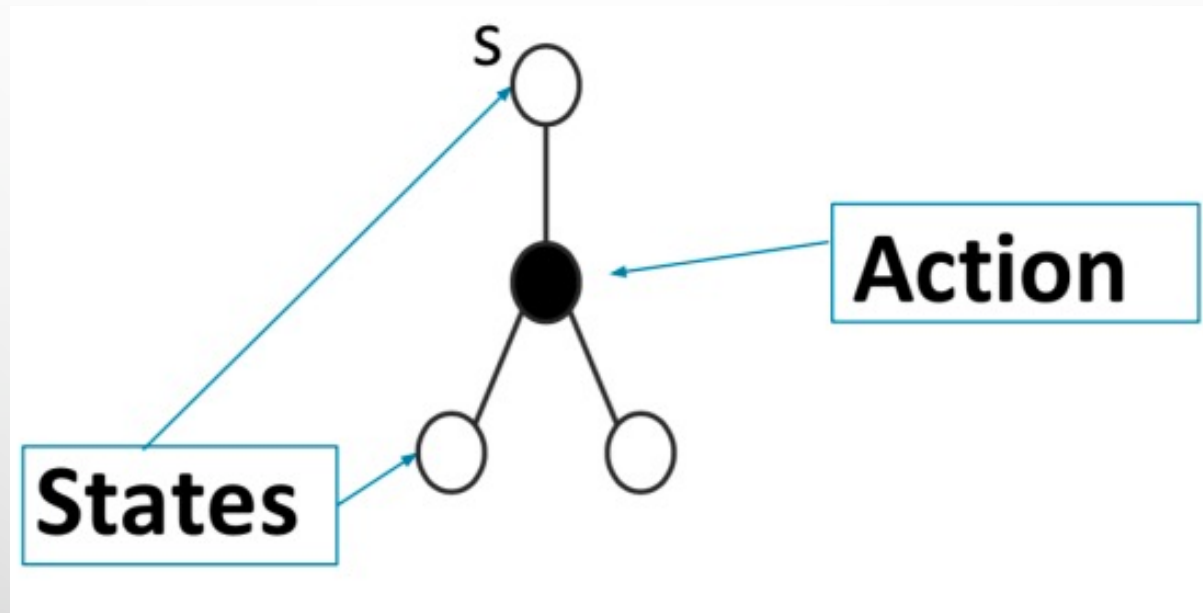
- Define $G_{i,t} = r_{i,t} + \gamma r_{i,t+1} + \gamma^2 r_{i,t+2} + \dots$ as return from time step t onwards in i th episode
- For state s visited at time step t in episode i
 - Increment counter of total visits: $N(s) = N(s) + 1$
 - Update estimate

$$V^\pi(s) = V^\pi(s) \frac{N(s) - 1}{N(s)} + \frac{G_{i,t}}{N(s)} = V^\pi(s) + \frac{1}{N(s)} (G_{i,t} - V^\pi(s))$$

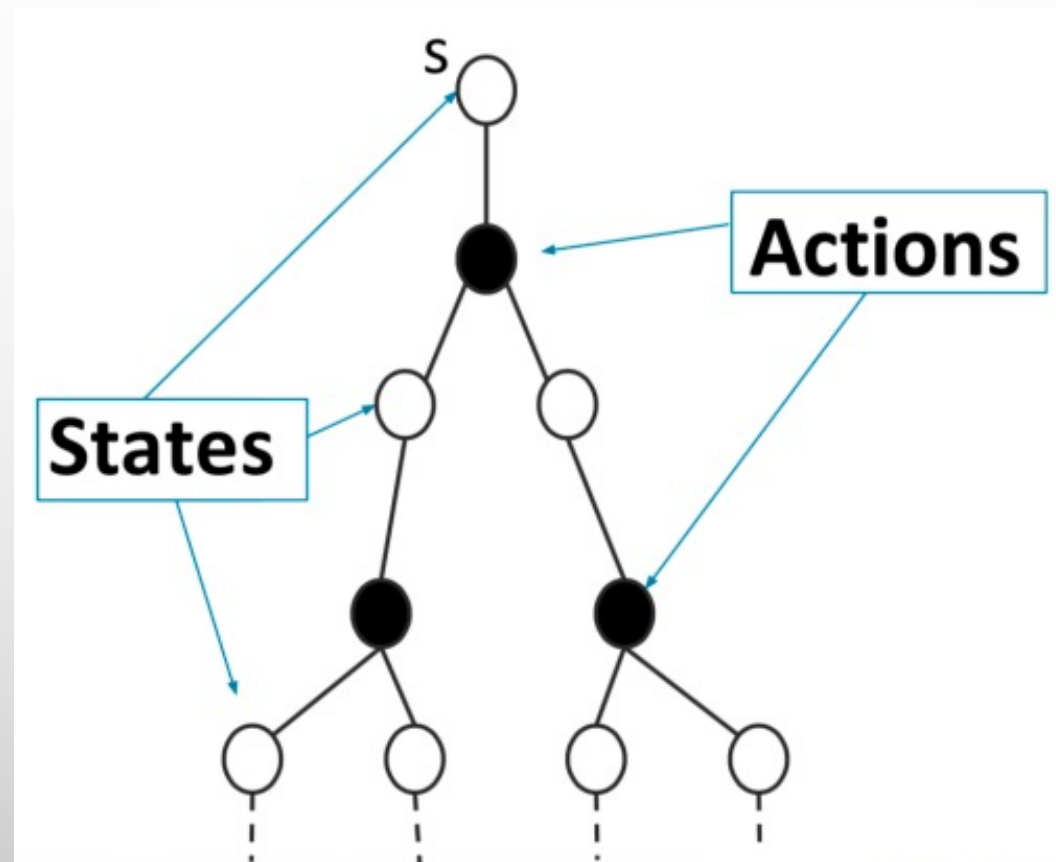
Policy Evaluation Diagram



Policy Evaluation Diagram

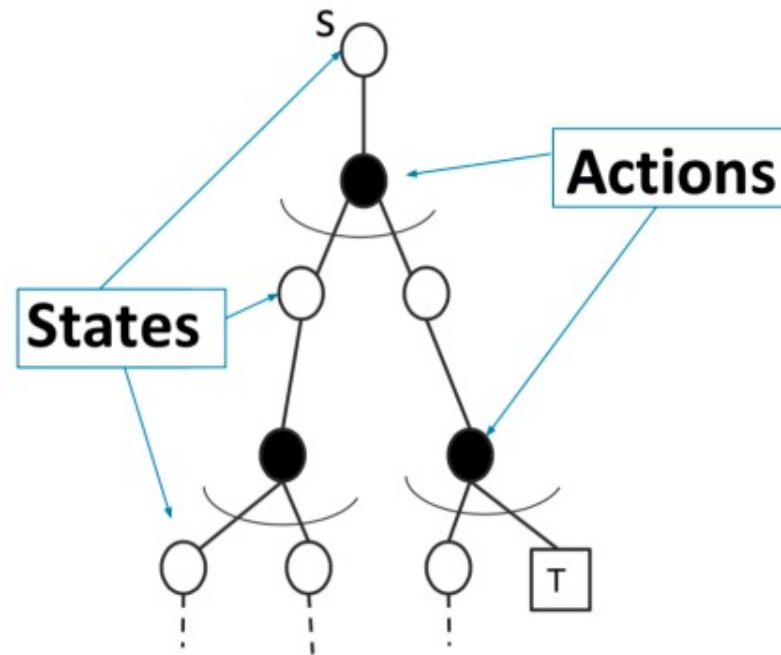


Policy Evaluation Diagram



Policy Evaluation Diagram

$$V^\pi(s) = V^\pi(s) + \alpha(G_{i,t} - V^\pi(s))$$

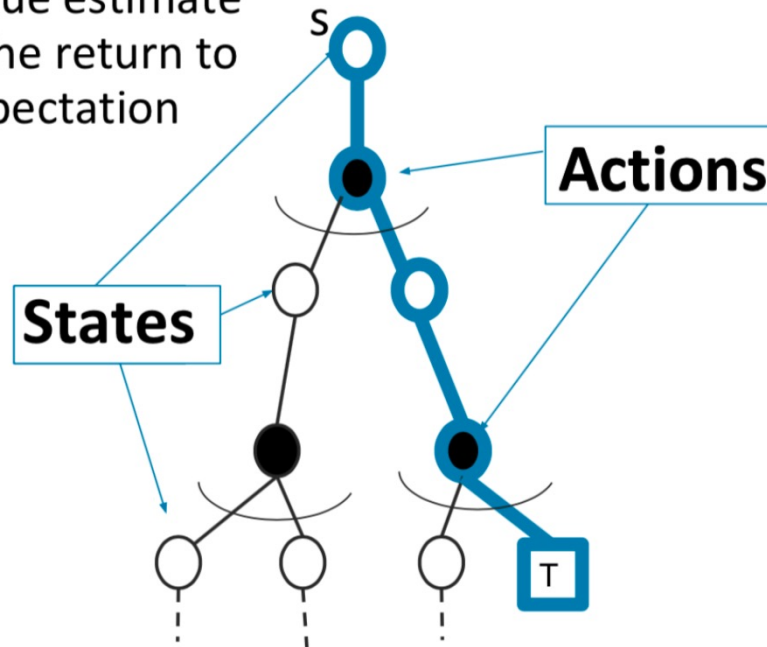


 = Expectation
 = **Terminal state**

Policy Evaluation Diagram

$$V^\pi(s) = V^\pi(s) + \alpha(G_{i,t} - V^\pi(s))$$

MC updates the value estimate using a **sample** of the return to approximate an expectation



 = Expectation
 = Terminal state