

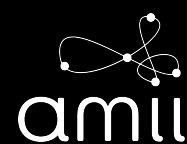
ESSENTIALS OF RL

Reinforcement Learning: Lecture 2

3rd Nepal Winter School in AI

24th Dec 2021

Abhishek Naik



UNIVERSITY OF
ALBERTA



OUTLINE

- ▶ Dynamic Programming (DP)
- ▶ Temporal-Difference (TD) Learning
- ▶ Model-based RL
- ▶ Policy Optimization

REINFORCEMENT LEARNING

REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal

REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ *via trial and error*

REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ *via trial and error*
 - ▶ *with potentially delayed rewards.*

REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ via *trial and error*
 - ▶ with potentially *delayed rewards*.

$$\pi_0 \longrightarrow \pi^*$$

REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ via *trial and error*
 - ▶ with potentially *delayed rewards*.

π_0

π^*

REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ via *trial and error*
 - ▶ with potentially *delayed rewards*.

π_0

π_1

π^*

REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ via *trial and error*
 - ▶ with potentially *delayed rewards*.

π_0

π_1

π_2

π^*

REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ via *trial and error*
 - ▶ with potentially *delayed rewards*.

$$\pi_0 \quad \pi_1 \quad \pi_2 \quad \pi_{n-1} \quad \pi^*$$

REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ via *trial and error*
 - ▶ with potentially *delayed rewards*.



REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ via *trial and error*
 - ▶ with potentially *delayed rewards*.



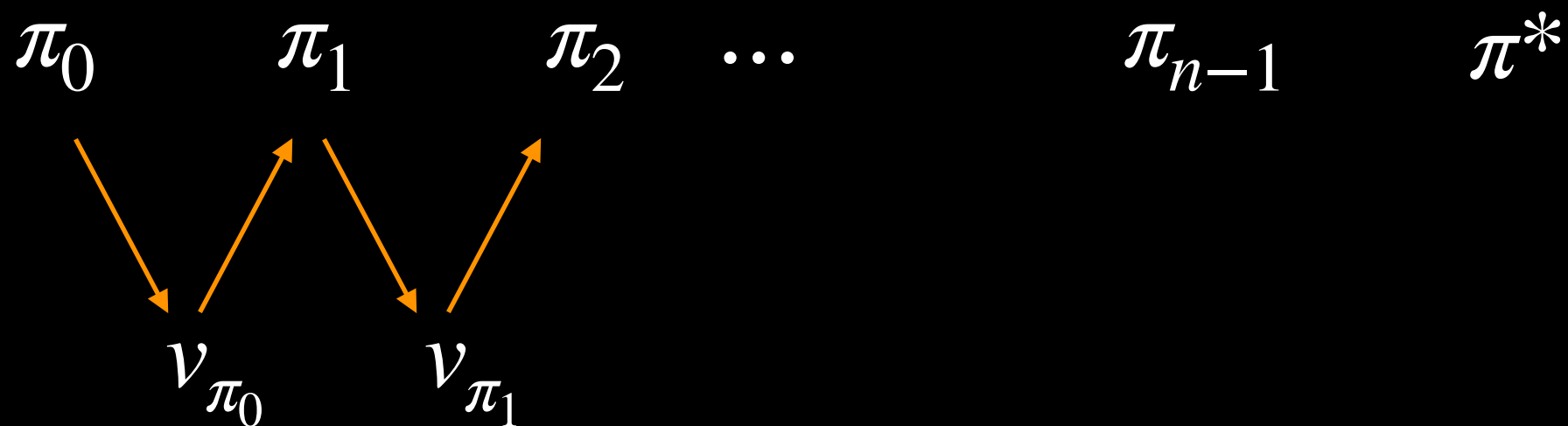
REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ via *trial and error*
 - ▶ with potentially *delayed rewards*.



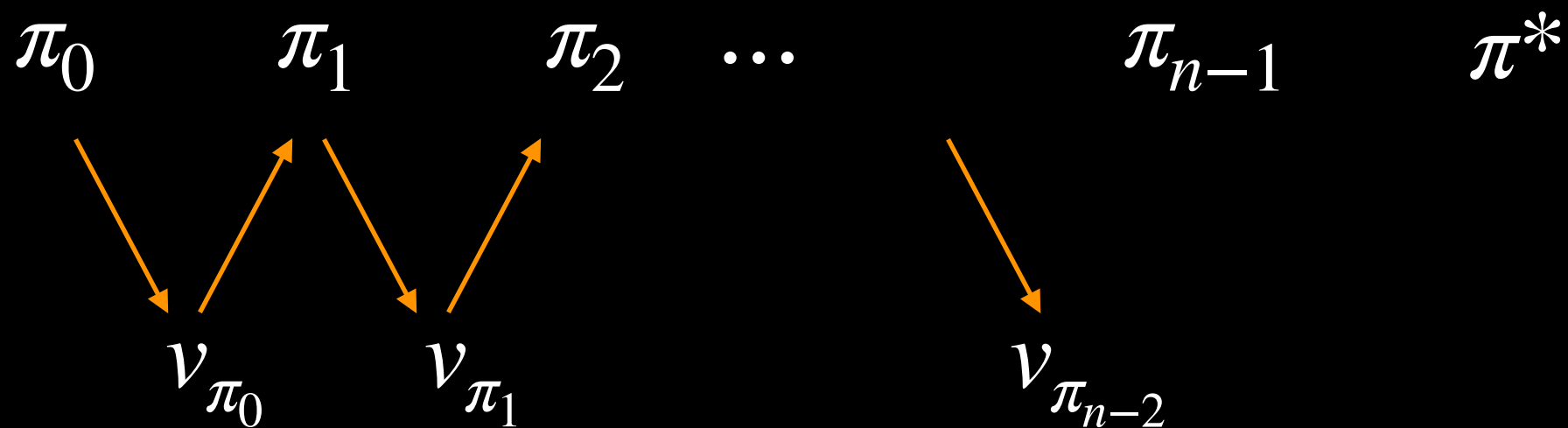
REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ via *trial and error*
 - ▶ with potentially *delayed rewards*.



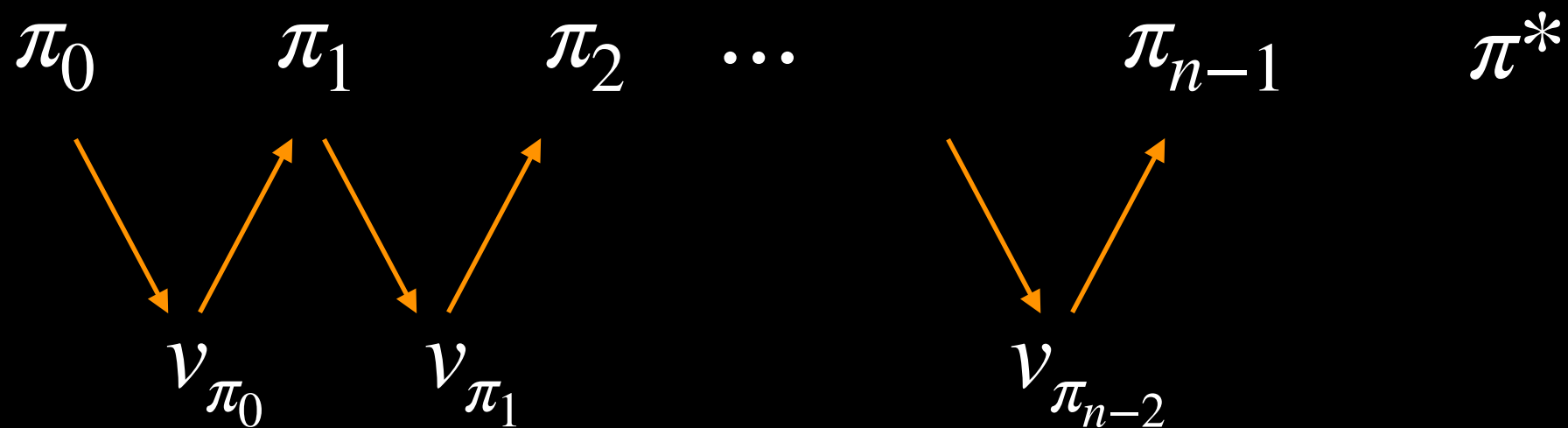
REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ via *trial and error*
 - ▶ with potentially *delayed rewards*.



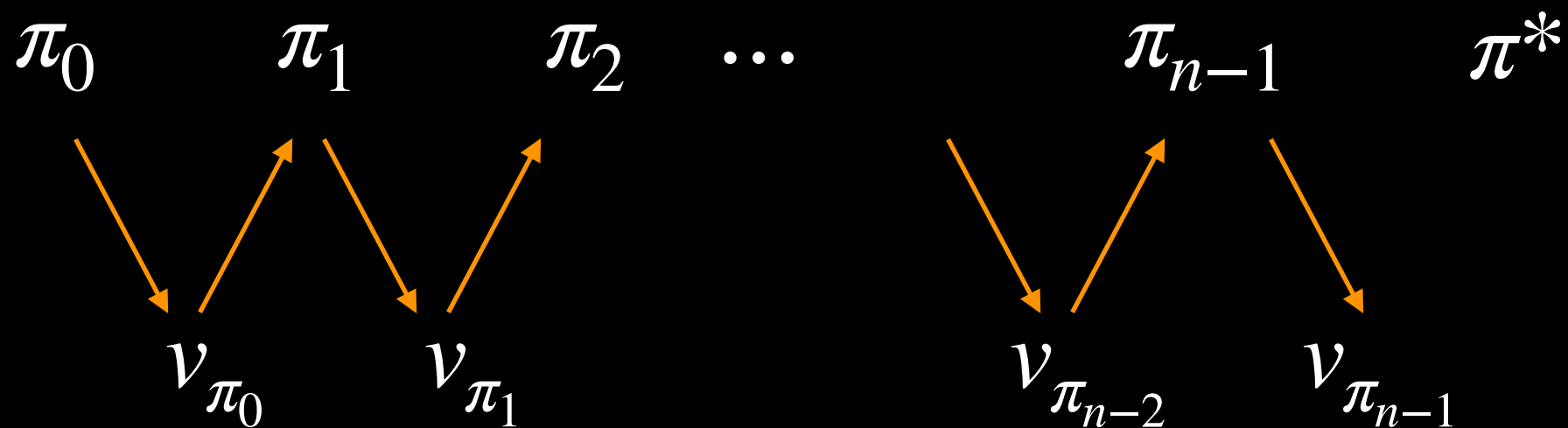
REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ via *trial and error*
 - ▶ with potentially *delayed rewards*.



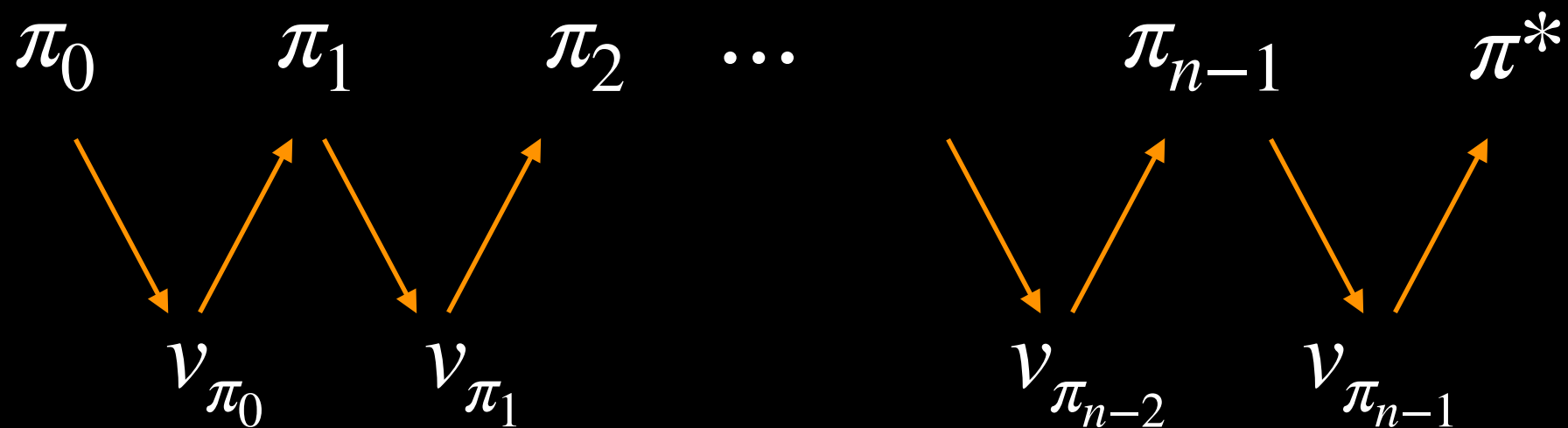
REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ via *trial and error*
 - ▶ with potentially *delayed rewards*.



REINFORCEMENT LEARNING

- ▶ Goal: learning some behaviour to maximize a numerical reward signal
 - ▶ via *trial and error*
 - ▶ with potentially *delayed rewards*.



REINFORCEMENT LEARNING

REINFORCEMENT LEARNING

- ▶ Sub-goal: evaluate a policy

REINFORCEMENT LEARNING

- ▶ Sub-goal: evaluate a policy
 - ▶ estimate the value function for a given policy

REINFORCEMENT LEARNING

- ▶ Sub-goal: evaluate a policy
 - ▶ estimate the value function for a given policy

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s]$$

REINFORCEMENT LEARNING

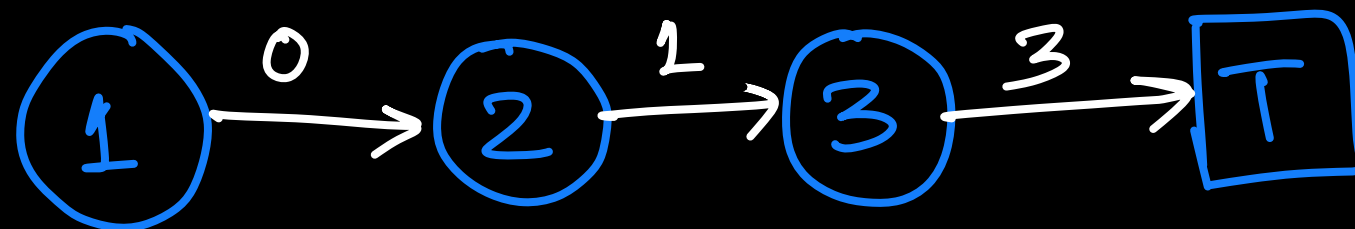
- ▶ Sub-goal: evaluate a policy
 - ▶ estimate the value function for a given policy

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] \qquad G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

REINFORCEMENT LEARNING

- ▶ Sub-goal: evaluate a policy
 - ▶ estimate the value function for a given policy

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] \qquad G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

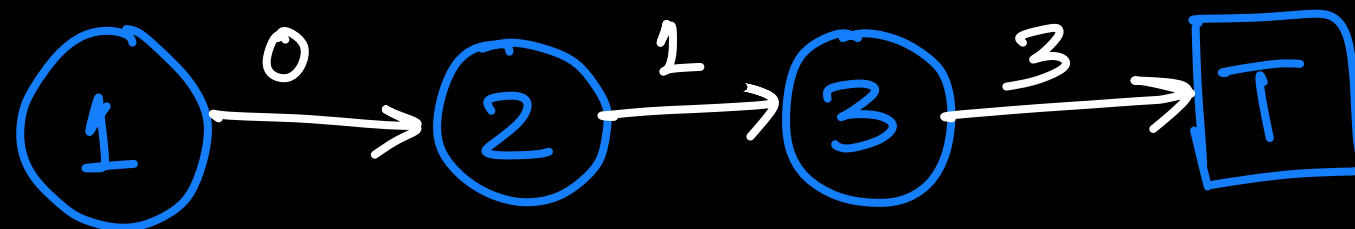


REINFORCEMENT LEARNING

- ▶ Sub-goal: evaluate a policy
 - ▶ estimate the value function for a given policy

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] \qquad G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

$$= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi}(s')]$$



FUNDAMENTALS

$$v_{\pi}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + v_{\pi}(s')]$$

FUNDAMENTALS

$$v_{\pi}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + v_{\pi}(s')]$$

Bellman equation

FUNDAMENTALS

$$v_{\pi}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + v_{\pi}(s')]$$

Bellman equation

Evaluating the optimal policy, π^* :

FUNDAMENTALS

$$v_{\pi}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + v_{\pi}(s')]$$

Bellman equation

Evaluating the optimal policy, π^* :

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

FUNDAMENTALS

$$v_{\pi}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + v_{\pi}(s')]$$

Bellman equation

Evaluating the optimal policy, π^* :

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

Bellman optimality equation

FUNDAMENTALS

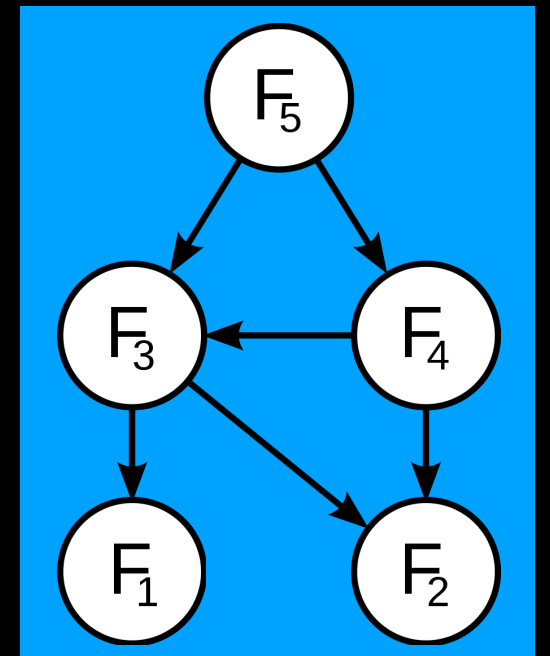
$$v_{\pi}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + v_{\pi}(s')]$$

Bellman equation

Evaluating the optimal policy, π^* :

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

Bellman optimality equation



FUNDAMENTALS

$$\begin{aligned} q_{\pi}(s, a) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi}(s')] \\ &= \sum_{s', r} p(s', r \mid s, a) \left[r + \sum_{a'} \pi(a' \mid s') q_{\pi}(s', a') \right] \end{aligned}$$

$$\begin{aligned} v_{\pi}(s) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] \\ &= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi}(s')] \\ v_{\pi}(s) &= \sum_a \pi(a \mid s) q_{\pi}(s, a) \end{aligned}$$

$$v_{\pi^*}(s) = \max_a \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi^*}(s')]$$

$$q_{\pi^*}(s, a) = \sum_{s', r} p(s', r \mid s, a) [r + \max_{a'} q_{\pi^*}(s', a')]$$

FUNDAMENTALS

$$\begin{aligned} q_{\pi}(s, a) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi}(s')] \\ &= \sum_{s', r} p(s', r \mid s, a) \left[r + \sum_{a'} \pi(a' \mid s') q_{\pi}(s', a') \right] \end{aligned}$$

$$\begin{aligned} v_{\pi}(s) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] \\ &= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi}(s')] \\ v_{\pi}(s) &= \sum_a \pi(a \mid s) q_{\pi}(s, a) \end{aligned}$$

$$v_{\pi^*}(s) = \max_a \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi^*}(s')]$$

$$q_{\pi^*}(s, a) = \sum_{s', r} p(s', r \mid s, a) [r + \max_{a'} q_{\pi^*}(s', a')]$$

FUNDAMENTALS

$$\begin{aligned} q_{\pi}(s, a) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi}(s')] \end{aligned}$$

$$= \sum_{s', r} p(s', r \mid s, a) \left[r + \sum_{a'} \pi(a' \mid s') q_{\pi}(s', a') \right]$$

$$\begin{aligned} v_{\pi}(s) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] \\ &= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi}(s')] \end{aligned}$$

$$v_{\pi}(s) = \sum_a \pi(a \mid s) q_{\pi}(s, a)$$

$$v_{\pi^*}(s) = \max_a \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi^*}(s')]$$

$$q_{\pi^*}(s, a) = \sum_{s', r} p(s', r \mid s, a) [r + \max_{a'} q_{\pi^*}(s', a')]$$

FUNDAMENTALS

$$\begin{aligned} q_{\pi}(s, a) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi}(s')] \\ &= \sum_{s', r} p(s', r \mid s, a) \left[r + \sum_{a'} \pi(a' \mid s') q_{\pi}(s', a') \right] \end{aligned}$$

$$\begin{aligned} v_{\pi}(s) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] \\ &= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi}(s')] \\ v_{\pi}(s) &= \sum_a \pi(a \mid s) q_{\pi}(s, a) \end{aligned}$$

$$v_{\pi^*}(s) = \max_a \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi^*}(s')]$$

$$q_{\pi^*}(s, a) = \sum_{s', r} p(s', r \mid s, a) [r + \max_{a'} q_{\pi^*}(s', a')]$$

FUNDAMENTALS

$$\begin{aligned} q_{\pi}(s, a) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi}(s')] \\ &= \sum_{s', r} p(s', r \mid s, a) \left[r + \sum_{a'} \pi(a' \mid s') q_{\pi}(s', a') \right] \end{aligned}$$

$$\begin{aligned} v_{\pi}(s) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] \\ &= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi}(s')] \\ v_{\pi}(s) &= \sum_a \pi(a \mid s) q_{\pi}(s, a) \end{aligned}$$

$$v_{\pi^*}(s) = \max_a \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi^*}(s')]$$

$$q_{\pi^*}(s, a) = \sum_{s', r} p(s', r \mid s, a) [r + \max_{a'} q_{\pi^*}(s', a')]$$

FUNDAMENTALS

$$\begin{aligned} q_{\pi}(s, a) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] \\ &= \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi}(s')] \\ &= \sum_{s', r} p(s', r \mid s, a) \left[r + \sum_{a'} \pi(a' \mid s') q_{\pi}(s', a') \right] \end{aligned}$$

$$\begin{aligned} v_{\pi}(s) &\doteq \mathbb{E}_{\pi}[G_t \mid S_t = s] \\ &= \sum_a \pi(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi}(s')] \\ v_{\pi}(s) &= \sum_a \pi(a \mid s) q_{\pi}(s, a) \end{aligned}$$

$$v_{\pi^*}(s) = \max_a \sum_{s', r} p(s', r \mid s, a) [r + v_{\pi^*}(s')]$$

$$q_{\pi^*}(s, a) = \sum_{s', r} p(s', r \mid s, a) [r + \max_{a'} q_{\pi^*}(s', a')]$$

CORE IDEAS OF DYNAMIC PROGRAMMING (DP)

CORE IDEAS OF DYNAMIC PROGRAMMING (DP)

► Goal:

CORE IDEAS OF DYNAMIC PROGRAMMING (DP)

- ▶ Goal:

- ▶ evaluate a given policy (the *prediction* problem), or

CORE IDEAS OF DYNAMIC PROGRAMMING (DP)

- ▶ Goal:

- ▶ evaluate a given policy (the *prediction* problem), or
- ▶ find the best policy (the *control* problem).

CORE IDEAS OF DYNAMIC PROGRAMMING (DP)

- ▶ Goal:

- ▶ evaluate a given policy (the *prediction* problem), or
- ▶ find the best policy (the *control* problem).

- ▶ Methodology:

CORE IDEAS OF DYNAMIC PROGRAMMING (DP)

- ▶ **Goal:**

- ▶ evaluate a given policy (the *prediction* problem), or
- ▶ find the best policy (the *control* problem).

- ▶ **Methodology:**

- ▶ Given a *model* of the world,

CORE IDEAS OF DYNAMIC PROGRAMMING (DP)

- ▶ **Goal:**

- ▶ evaluate a given policy (the *prediction* problem), or
- ▶ find the best policy (the *control* problem).

- ▶ **Methodology:**

- ▶ Given a *model* of the world,
- ▶ *iteratively* improve estimates

CORE IDEAS OF DYNAMIC PROGRAMMING (DP)

▶ Goal:

- ▶ evaluate a given policy (the *prediction* problem), or
- ▶ find the best policy (the *control* problem).

▶ Methodology:

- ▶ Given a *model* of the world,
- ▶ *iteratively* improve estimates
- ▶ with the help of *bootstrapping*.

CORE IDEAS OF DYNAMIC PROGRAMMING (DP)

▶ Goal:

- ▶ evaluate a given policy (the *prediction* problem), or
- ▶ find the best policy (the *control* problem).

▶ Methodology:

- ▶ Given a *model* of the world,
- ▶ *iteratively* improve estimates
- ▶ with the help of *bootstrapping*.

$$v_{t+1}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + v_t(s')]$$

CORE IDEAS OF DYNAMIC PROGRAMMING (DP)

▶ Goal:

- ▶ evaluate a given policy (the *prediction* problem), or
- ▶ find the best policy (the *control* problem).

▶ Methodology:

- ▶ Given a model of the world,
- ▶ *iteratively* improve estimates
- ▶ with the help of *bootstrapping*.

$$v_{t+1}(s) = \sum_a \pi(a | s) \sum_{s', r} \underline{p(s', r | s, a)} [r + v_t(s')]$$

CORE IDEAS OF DYNAMIC PROGRAMMING (DP)

► Goal:

- evaluate a given policy (the *prediction* problem), or
- find the best policy (the *control* problem).

► Methodology:

- Given a *model* of the world,
- *iteratively* improve estimates
- with the help of *bootstrapping*.

$$v_{t+1}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + v_t(s')]$$

CORE IDEAS OF DYNAMIC PROGRAMMING (DP)

► Goal:

- evaluate a given policy (the *prediction* problem), or
- find the best policy (the *control* problem).

► Methodology:

- Given a *model* of the world,
- *iteratively* improve estimates
- with the help of *bootstrapping*.

$$v_{t+1}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + v_t(s')]$$

WE DO NOT USUALLY HAVE A MODEL OF THE WORLD

WE DO NOT USUALLY HAVE A MODEL OF THE WORLD

- ▶ But we can (and do!) interact with the world.

WE DO NOT USUALLY HAVE A MODEL OF THE WORLD

- ▶ But we can (and do!) interact with the world.
- ▶ Can we use that experience to improve our estimates?

WE DO NOT USUALLY HAVE A MODEL OF THE WORLD

- ▶ But we can (and do!) interact with the world.
- ▶ Can we use that experience to improve our estimates?

$$v_{t+1}(s) \rightarrow [r + v_t(s')]$$

WE DO NOT USUALLY HAVE A MODEL OF THE WORLD

- ▶ But we can (and do!) interact with the world.
- ▶ Can we use that experience to improve our estimates?

$$v_{t+1}(s) \rightarrow [r + v_t(s')]$$

How can we improve estimates
from a *stream* of data?

WE DO NOT USUALLY HAVE A MODEL OF THE WORLD

- ▶ But we can (and do!) interact with the world.
- ▶ Can we use that experience to improve our estimates?

$$v_{t+1}(s) \rightarrow [r + v_t(s')]$$

How can we improve estimates
from a *stream* of data?

$$S_0, A_0, R_1, S_1, \dots, S_t, A_t, R_{t+1}, S_{t+1}, \dots$$

IMPROVING ESTIMATES FROM A STREAM OF DATA

IMPROVING ESTIMATES FROM A STREAM OF DATA

- ▶ An example for intuition:
 - ▶ Compute the average of a stream of samples

IMPROVING ESTIMATES FROM A STREAM OF DATA

- ▶ An example for intuition:
 - ▶ Compute the average of a stream of samples

$$x_1, x_2, x_3, \dots, x_N, x_{N+1}, \dots$$

IMPROVING ESTIMATES FROM A STREAM OF DATA

- ▶ An example for intuition:

- ▶ Compute the average of a stream of samples

$$\bar{x}_{N+1} = \frac{1}{N+1} \left(\sum_{i=1}^N x_i + x_{N+1} \right)$$

$$x_1, x_2, x_3, \dots, x_N, x_{N+1}, \dots$$

IMPROVING ESTIMATES FROM A STREAM OF DATA

- ▶ An example for intuition:
 - ▶ Compute the average of a stream of samples

$x_1, x_2, x_3, \dots, x_N, x_{N+1}, \dots$

$$\begin{aligned}\bar{x}_{N+1} &= \frac{1}{N+1} \left(\sum_{i=1}^N x_i + x_{N+1} \right) \\ &= \frac{N}{N+1} \frac{1}{N} \sum_{i=1}^N x_i + \frac{1}{N+1} x_{N+1}\end{aligned}$$

IMPROVING ESTIMATES FROM A STREAM OF DATA

- ▶ An example for intuition:
 - ▶ Compute the average of a stream of samples

$x_1, x_2, x_3, \dots, x_N, x_{N+1}, \dots$

$$\begin{aligned}\bar{x}_{N+1} &= \frac{1}{N+1} \left(\sum_{i=1}^N x_i + x_{N+1} \right) \\ &= \frac{N}{N+1} \frac{1}{N} \sum_{i=1}^N x_i + \frac{1}{N+1} x_{N+1} \\ &= \frac{N}{N+1} \bar{x}_N + \frac{1}{N+1} x_{N+1}\end{aligned}$$

IMPROVING ESTIMATES FROM A STREAM OF DATA

- ▶ An example for intuition:
 - ▶ Compute the average of a stream of samples

$x_1, x_2, x_3, \dots, x_N, x_{N+1}, \dots$

$$\begin{aligned}\bar{x}_{N+1} &= \frac{1}{N+1} \left(\sum_{i=1}^N x_i + x_{N+1} \right) \\ &= \frac{N}{N+1} \frac{1}{N} \sum_{i=1}^N x_i + \frac{1}{N+1} x_{N+1} \\ &= \frac{N}{N+1} \bar{x}_N + \frac{1}{N+1} x_{N+1} \\ &= \left(1 - \frac{1}{N+1} \right) \bar{x}_N + \frac{1}{N+1} x_{N+1}\end{aligned}$$

IMPROVING ESTIMATES FROM A STREAM OF DATA

- ▶ An example for intuition:
 - ▶ Compute the average of a stream of samples

$x_1, x_2, x_3, \dots, x_N, x_{N+1}, \dots$

$$\begin{aligned}\bar{x}_{N+1} &= \frac{1}{N+1} \left(\sum_{i=1}^N x_i + x_{N+1} \right) \\ &= \frac{N}{N+1} \frac{1}{N} \sum_{i=1}^N x_i + \frac{1}{N+1} x_{N+1} \\ &= \frac{N}{N+1} \bar{x}_N + \frac{1}{N+1} x_{N+1} \\ &= \left(1 - \frac{1}{N+1} \right) \bar{x}_N + \frac{1}{N+1} x_{N+1} \\ &= \bar{x}_N + \frac{1}{N+1} (x_{N+1} - \bar{x}_N)\end{aligned}$$

IMPROVING ESTIMATES FROM A STREAM OF DATA

- ▶ An example for intuition:
 - ▶ Compute the average of a stream of samples

$x_1, x_2, x_3, \dots, x_N, x_{N+1}, \dots$

$$\begin{aligned}\bar{x}_{N+1} &= \frac{1}{N+1} \left(\sum_{i=1}^N x_i + x_{N+1} \right) \\&= \frac{N}{N+1} \frac{1}{N} \sum_{i=1}^N x_i + \frac{1}{N+1} x_{N+1} \\&= \frac{N}{N+1} \bar{x}_N + \frac{1}{N+1} x_{N+1} \\&= \left(1 - \frac{1}{N+1} \right) \bar{x}_N + \frac{1}{N+1} x_{N+1} \\&= \bar{x}_N + \frac{1}{N+1} (x_{N+1} - \bar{x}_N) \\ \bar{x}_{N+1} &= \bar{x}_N + \alpha (x_{N+1} - \bar{x}_N)\end{aligned}$$

IMPROVING ESTIMATES FROM A STREAM OF DATA

$$\bar{x}_{N+1} = \bar{x}_N + \alpha (x_{N+1} - \bar{x}_N)$$

IMPROVING ESTIMATES FROM A STREAM OF DATA

$$\bar{x}_{N+1} = \bar{x}_N + \alpha (x_{N+1} - \bar{x}_N)$$

```
new_estimate = old_estimate + stepsize*(new_target - old_estimate)
```

IMPROVING ESTIMATES FROM A STREAM OF DATA

$$\bar{x}_{N+1} = \bar{x}_N + \alpha (x_{N+1} - \bar{x}_N)$$

`new_estimate = old_estimate + stepsize * (new_target - old_estimate)`

DP
$$v_{t+1}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + v_t(s')]$$

IMPROVING ESTIMATES FROM A STREAM OF DATA

$$\bar{x}_{N+1} = \bar{x}_N + \alpha (x_{N+1} - \bar{x}_N)$$

`new_estimate = old_estimate + stepsize * (new_target - old_estimate)`

DP
$$v_{t+1}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + v_t(s')]$$

TD
$$v_{t+1}(s) = v_t(s) + \alpha \left[(r + v_t(s')) - v_t(s) \right]$$

TEMPORAL-DIFFERENCE (TD) LEARNING

TEMPORAL-DIFFERENCE (TD) LEARNING

Algorithm : Tabular TD learning to estimate v_π

Input: The target policy π

Algorithm parameters: step size $\alpha \in (0, 1]$

- 1 Initialize $V(s)$, for all $s \in \mathcal{S}$, arbitrarily (e.g., to zero)
 - 2 Observe initial state S
 - 3 **for** *each time step* **do**
 - 4 $A \leftarrow$ action according to π in S
 - 5 Take action A , observe R, S'
 - 6 $V(S) \leftarrow V(S) + \alpha [R + V(S') - V(S)]$
 - 7 $S \leftarrow S'$
 - 8 **end**
 - 9 return V
-

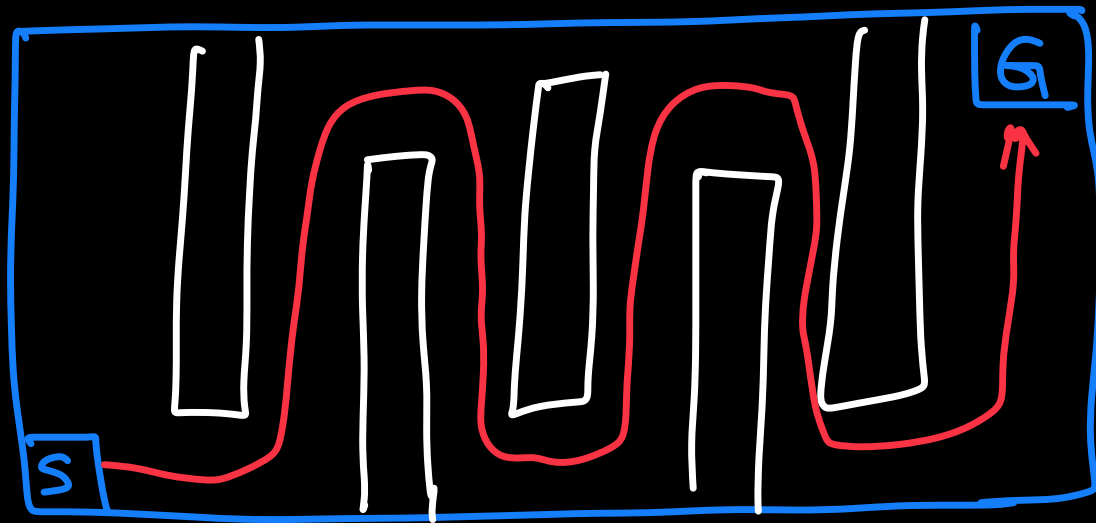
CONTROL: EXPLORATION VS EXPLOITATION

CONTROL: EXPLORATION VS EXPLOITATION

- ▶ Simple heuristic:
 - ▶ with a small probability, pick a random action

CONTROL: EXPLORATION VS EXPLOITATION

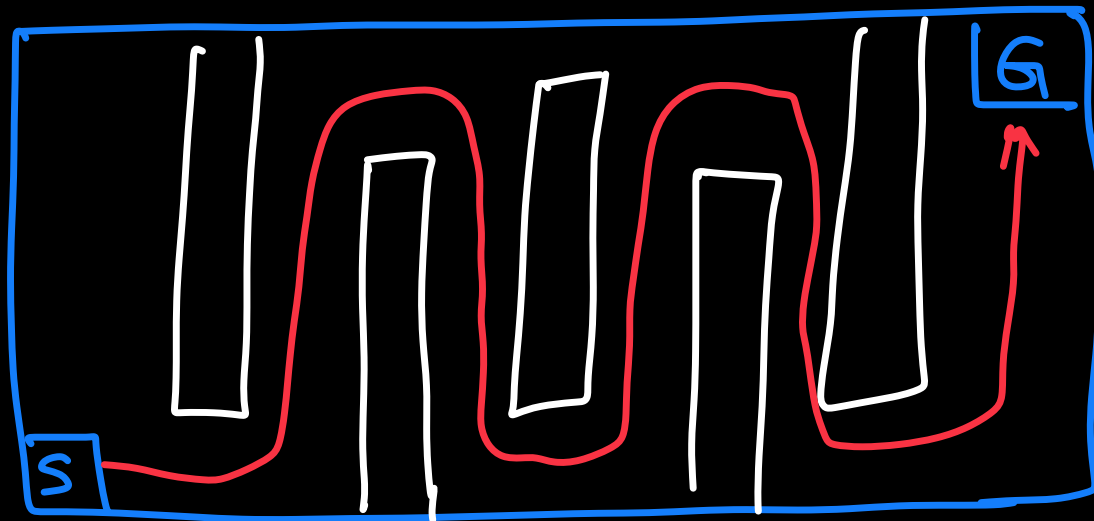
- ▶ Simple heuristic:
 - ▶ with a small probability, pick a random action



CONTROL: EXPLORATION VS EXPLOITATION

- ▶ Simple heuristic:

- ▶ with a small probability, pick a random action

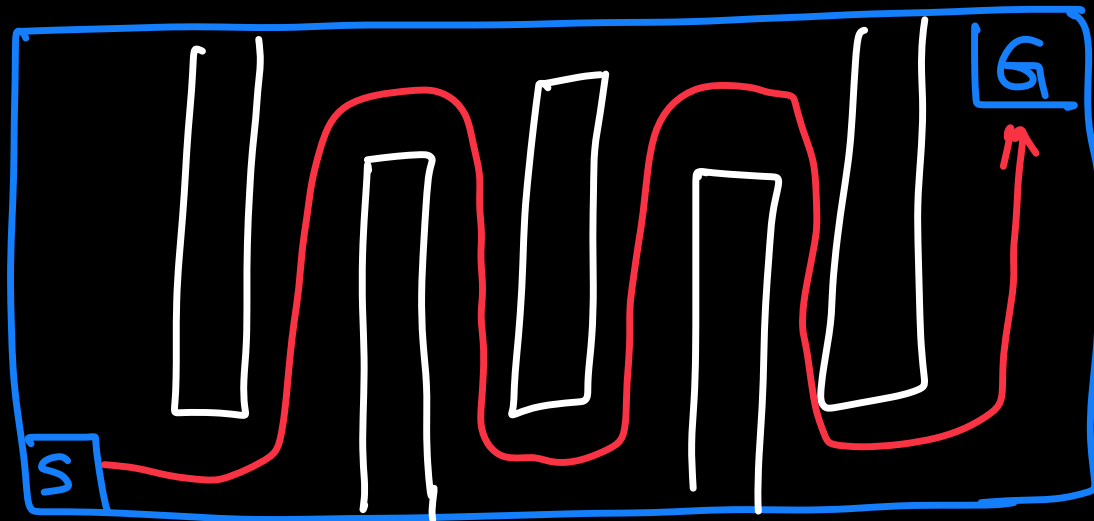


- With probability ϵ , pick an action randomly
- With probability $1-\epsilon$, pick the 'greedy' action

CONTROL: EXPLORATION VS EXPLOITATION

- ▶ Simple heuristic:

- ▶ with a small probability, pick a random action



- With probability ϵ , pick an action randomly
- With probability $1-\epsilon$, pick the 'greedy' action

ϵ -greedy action selection

CONTROL ALGORITHM: SARSA

Algorithm : SARSA to estimate $Q \approx Q_{\pi^*}$

Parameters: step size $\alpha \in (0, 1]$

- 1 Initialize $Q(s, a)$, for all $s \in \mathcal{S}, a \in \mathcal{A}$, arbitrarily (e.g., to zero)
 - 2 Observe initial state S
 - 3 **for** *each time step* **do**
 - 4 $A \leftarrow$ action in S according to policy derived from Q (e.g., ϵ -greedy)
 - 5 Take action A , observe R, S'
 - 6 $A' \leftarrow$ action in S' according to policy derived from Q (e.g., ϵ -greedy)
 - 7 $Q(S, A) \leftarrow Q(S, A) + \alpha [R + Q(S', A') - Q(S, A)]$
 - 8 $S \leftarrow S'$
 - 9 **end**
-

CONTROL ALGORITHM: SARSA

Algorithm : SARSA to estimate $Q \approx Q_{\pi^*}$

Parameters: step size $\alpha \in (0, 1]$

- 1 Initialize $Q(s, a)$, for all $s \in \mathcal{S}, a \in \mathcal{A}$, arbitrarily (e.g., to zero)
 - 2 Observe initial state S
 - 3 **for** *each time step* **do**
 - 4 $A \leftarrow$ action in S according to policy derived from Q (e.g., ϵ -greedy)
 - 5 Take action A , observe R, S'
 - 6 $A' \leftarrow$ action in S' according to policy derived from Q (e.g., ϵ -greedy)
 - 7 $Q(S, A) \leftarrow Q(S, A) + \alpha [R + Q(S', A') - Q(S, A)]$
 - 8 $S \leftarrow S'$
 - 9 **end**
-

“On-policy”

CONTROL ALGORITHM: Q-LEARNING

Algorithm : Q-learning to estimate $Q \approx Q_{\pi^*}$

Parameters: step size $\alpha \in (0, 1]$

- 1 Initialize $Q(s, a)$, for all $s \in \mathcal{S}, a \in \mathcal{A}$, arbitrarily (e.g., to zero)
 - 2 Observe initial state S
 - 3 **for** *each time step* **do**
 - 4 $A \leftarrow$ action in S according to policy derived from Q (e.g., ϵ -greedy)
 - 5 Take action A , observe R, S'
 - 6 $A' \leftarrow$ action in S' according to the greedy policy derived from Q
 - 7 $Q(S, A) \leftarrow Q(S, A) + \alpha [R + Q(S', A') - Q(S, A)]$
 - 8 $S \leftarrow S'$
 - 9 **end**
-

CONTROL ALGORITHM: Q-LEARNING

Algorithm : Q-learning to estimate $Q \approx Q_{\pi^*}$

Parameters: step size $\alpha \in (0, 1]$

- 1 Initialize $Q(s, a)$, for all $s \in \mathcal{S}, a \in \mathcal{A}$, arbitrarily (e.g., to zero)
 - 2 Observe initial state S
 - 3 **for** *each time step* **do**
 - 4 $A \leftarrow$ action in S according to policy derived from Q (e.g., ϵ -greedy)
 - 5 Take action A , observe R, S'
 - 6 $A' \leftarrow$ action in S' according to the greedy policy derived from Q
 - 7 $Q(S, A) \leftarrow Q(S, A) + \alpha [R + Q(S', A') - Q(S, A)]$
 - 8 $S \leftarrow S'$
 - 9 **end**
-


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

CONTROL ALGORITHM: Q-LEARNING

Algorithm : Q-learning to estimate $Q \approx Q_{\pi^*}$

Parameters: step size $\alpha \in (0, 1]$

- 1 Initialize $Q(s, a)$, for all $s \in \mathcal{S}, a \in \mathcal{A}$, arbitrarily (e.g., to zero)
 - 2 Observe initial state S
 - 3 **for** *each time step* **do**
 - 4 $A \leftarrow$ action in S according to policy derived from Q (e.g., ϵ -greedy)
 - 5 Take action A , observe R, S'
 - 6 $A' \leftarrow$ action in S' according to the greedy policy derived from Q
 - 7 $Q(S, A) \leftarrow Q(S, A) + \alpha [R + Q(S', A') - Q(S, A)]$
 - 8 $S \leftarrow S'$
 - 9 **end**
-

“Off-policy”


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

PLANNING USING A MODEL OF THE WORLD

PLANNING USING A MODEL OF THE WORLD

- ▶ We studied one way to use a model with DP:

PLANNING USING A MODEL OF THE WORLD

- ▶ We studied one way to use a model with DP:

$$v_{t+1}(s) = \sum_a \pi(a | s) \sum_{s', r} \underline{p(s', r | s, a)} [r + v_t(s')]$$

PLANNING USING A MODEL OF THE WORLD

- ▶ We studied one way to use a model with DP:

$$v_{t+1}(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + v_t(s')]$$

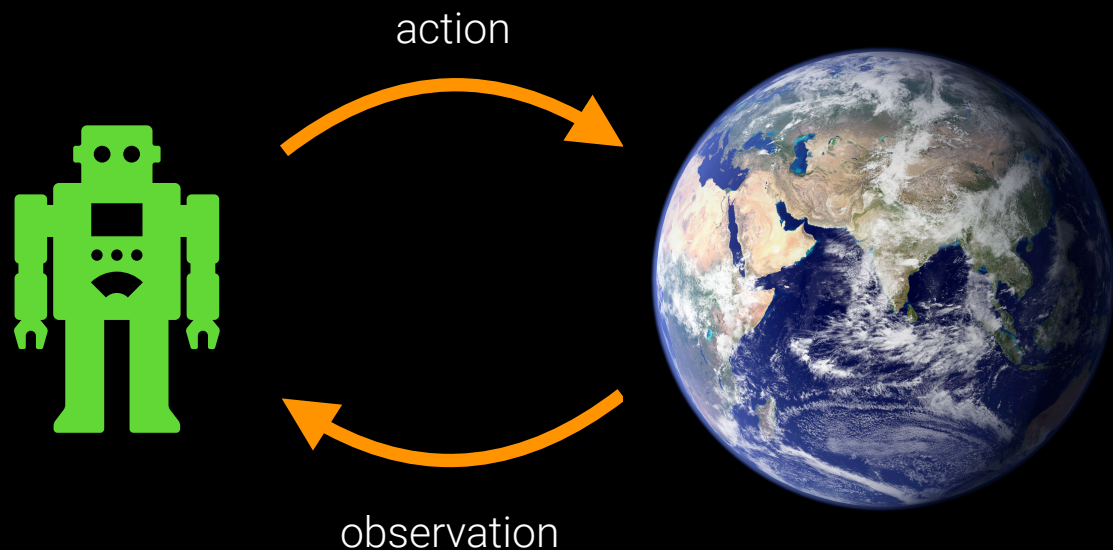
- ▶ We can also use the model as a substitute for real-world experience:

PLANNING USING A MODEL OF THE WORLD

- ▶ We studied one way to use a model with DP:

$$v_{t+1}(s) = \sum_a \pi(a | s) \sum_{s', r} \underline{p(s', r | s, a)} [r + v_t(s')]$$

- ▶ We can also use the model as a substitute for real-world experience:

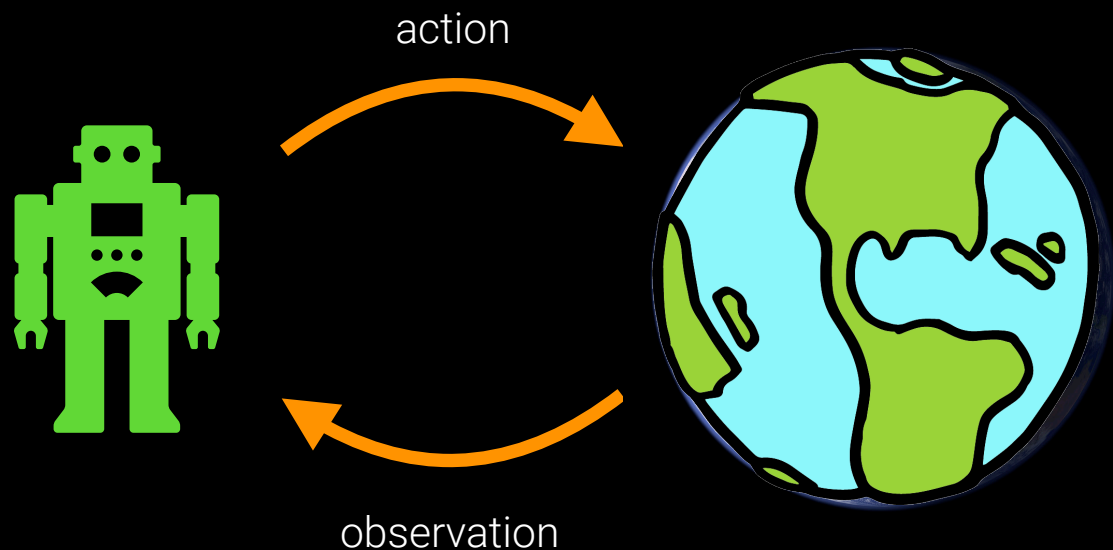


PLANNING USING A MODEL OF THE WORLD

- ▶ We studied one way to use a model with DP:

$$v_{t+1}(s) = \sum_a \pi(a | s) \sum_{s', r} \underline{p(s', r | s, a)} [r + v_t(s')]$$

- ▶ We can also use the model as a substitute for real-world experience:



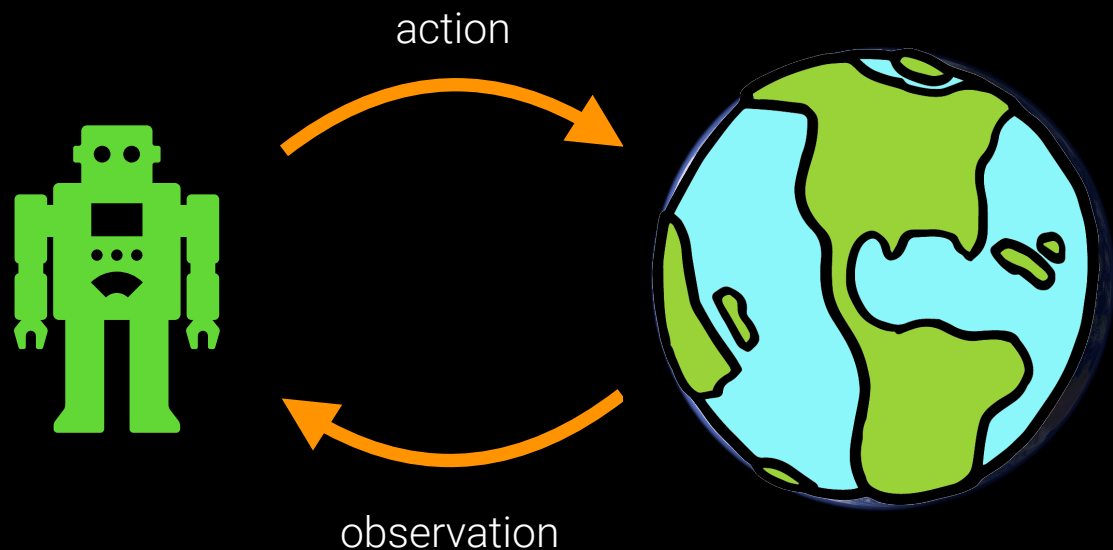
PLANNING USING A MODEL OF THE WORLD

- ▶ We studied one way to use a model with DP:

$$v_{t+1}(s) = \sum_a \pi(a | s) \sum_{s', r} \underline{p(s', r | s, a)} [r + v_t(s')]$$

- ▶ We can also use the model as a substitute for real-world experience:

$$v_{t+1}(S) = v_t(S) + \alpha [R + v_t(S') - v_t(S)]$$



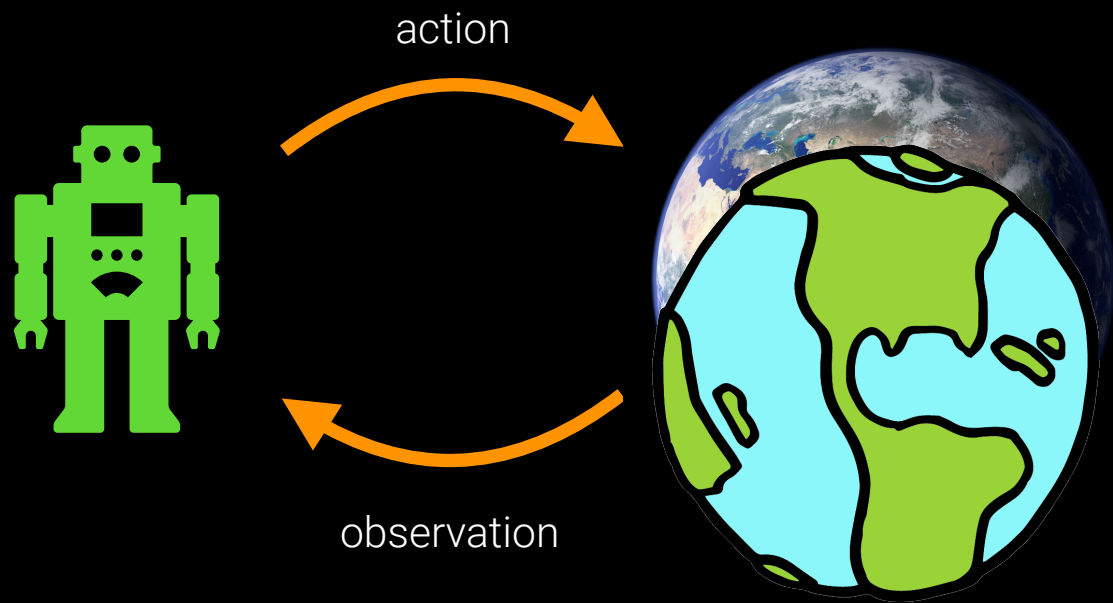
PLANNING USING A MODEL OF THE WORLD

- ▶ We studied one way to use a model with DP:

$$v_{t+1}(s) = \sum_a \pi(a | s) \sum_{s', r} \underline{p(s', r | s, a)} [r + v_t(s')]$$

- ▶ We can also use the model as a substitute for real-world experience:

$$v_{t+1}(S) = v_t(S) + \alpha [R + v_t(S') - v_t(S)]$$



LEARNING A MODEL FROM EXPERIENCE

LEARNING A MODEL FROM EXPERIENCE

- ▶ Where do models come from?

LEARNING A MODEL FROM EXPERIENCE

- ▶ Where do models come from?
 - ▶ They can also be learned from experience.

LEARNING A MODEL FROM EXPERIENCE

- ▶ Where do models come from?
 - ▶ They can also be learned from experience.

```
A a 1 B b 1 T
A b 0 C a 3 T
A a 1 B a 0 T
A b 0 C b 0 T
```

DYNA: INTEGRATING LEARNING AND PLANNING

{

Learning the model

{

Planning with
the learned model

DYNA: INTEGRATING LEARNING AND PLANNING

Algorithm : Dyna to estimate $Q \approx Q_{\pi^*}$

Parameters: step size $\alpha \in (0, 1]$

```
1 Initialize  $Q(s, a)$  and  $Model(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}$ 
2 Observe initial state  $S$ 
3 for each time step do
4      $A \leftarrow$  action in  $S$  according to policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
5     Take action  $A$ , observe  $R, S'$ 
6     { Update model:  $Model(S, A) \leftarrow S', R$                                 Learning the model
7     while time to plan do
8         {  $S_p \leftarrow$  random previously observed state                                Planning with
9         {  $A_p \leftarrow$  random previously picked action in  $S$                         the learned model
10        {  $R_p, S'_p \leftarrow Model(S_p, A_p)$ 
11        {  $Q(S_p, A_p) \leftarrow Q(S_p, A_p) + \alpha [R_p + \max_a Q(S'_p, a) - Q(S_p, A_p)]$ 
12        end
13         $S \leftarrow S'$ 
14 end
```

DYNA: INTEGRATING LEARNING AND PLANNING

Algorithm : Dyna to estimate $Q \approx Q_{\pi^*}$

Parameters: step size $\alpha \in (0, 1]$

```
1 Initialize  $Q(s, a)$  and  $Model(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}$ 
2 Observe initial state  $S$ 
3 for each time step do
4      $A \leftarrow$  action in  $S$  according to policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
5     Take action  $A$ , observe  $R, S'$ 
6     { Update model:  $Model(S, A) \leftarrow S', R$  Learning the model
7     while time to plan do
8         {  $S_p \leftarrow$  random previously observed state Planning with
9         {  $A_p \leftarrow$  random previously picked action in  $S$  the learned model
10        {  $R_p, S'_p \leftarrow Model(S_p, A_p)$ 
11        {  $Q(S_p, A_p) \leftarrow Q(S_p, A_p) + \alpha [R_p + \max_a Q(S'_p, a) - Q(S_p, A_p)]$ 
12        end
13         $S \leftarrow S'$ 
14 end
```

DYNA: INTEGRATING LEARNING AND PLANNING

Algorithm : Dyna to estimate $Q \approx Q_{\pi^*}$

Parameters: step size $\alpha \in (0, 1]$

```
1 Initialize  $Q(s, a)$  and  $Model(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}$ 
2 Observe initial state  $S$ 
3 for each time step do
4    $A \leftarrow$  action in  $S$  according to policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
5   Take action  $A$ , observe  $R, S'$ 
6   { Update model:  $Model(S, A) \leftarrow S', R$  Learning the model
7   while time to plan do
8     {  $S_p \leftarrow$  random previously observed state Planning with
9      $A_p \leftarrow$  random previously picked action in  $S$  the learned model
10     $R_p, S'_p \leftarrow Model(S_p, A_p)$ 
11     $Q(S_p, A_p) \leftarrow Q(S_p, A_p) + \alpha [R_p + \max_a Q(S'_p, a) - Q(S_p, A_p)]$ 
12  end
13   $S \leftarrow S'$ 
14 end
```

DYNA: INTEGRATING LEARNING AND PLANNING

Algorithm : Dyna to estimate $Q \approx Q_{\pi^*}$

Parameters: step size $\alpha \in (0, 1]$

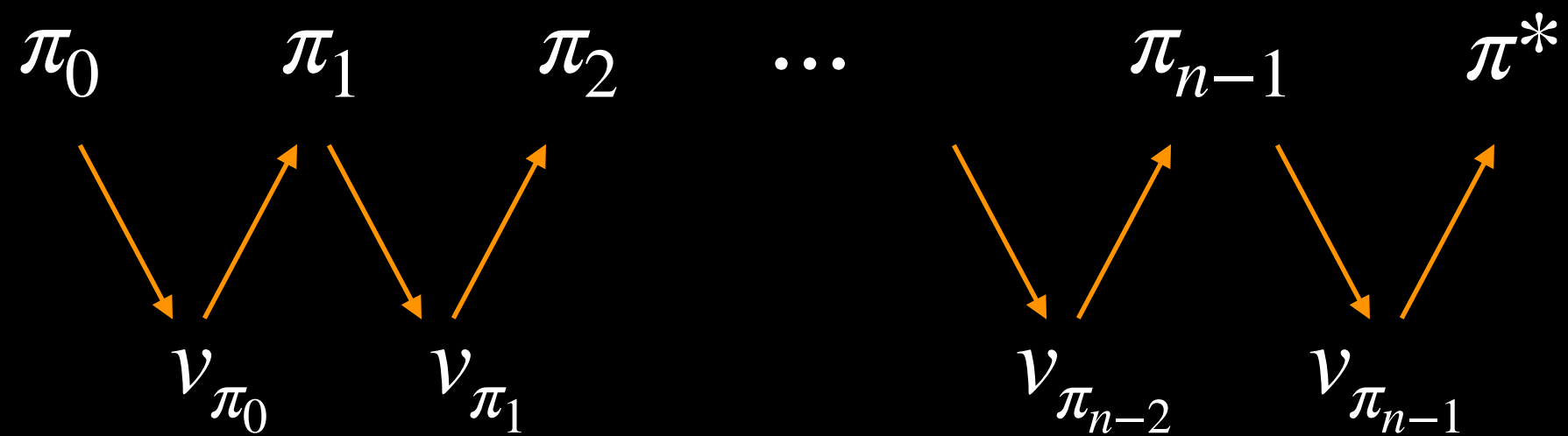
```
1 Initialize  $Q(s, a)$  and  $Model(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}$ 
2 Observe initial state  $S$ 
3 for each time step do
4    $A \leftarrow$  action in  $S$  according to policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
5   Take action  $A$ , observe  $R, S'$ 
6   { Update model:  $Model(S, A) \leftarrow S', R$                                 Learning the model
7   while time to plan do
8     {  $S_p \leftarrow$  random previously observed state                                Planning with
9     {  $A_p \leftarrow$  random previously picked action in  $S$                        the learned model
10    {  $R_p, S'_p \leftarrow Model(S_p, A_p)$ 
11    {  $Q(S_p, A_p) \leftarrow Q(S_p, A_p) + \alpha [R_p + \max_a Q(S'_p, a) - Q(S_p, A_p)]$ 
12    end
13     $S \leftarrow S'$ 
14 end
```

OUTLINE

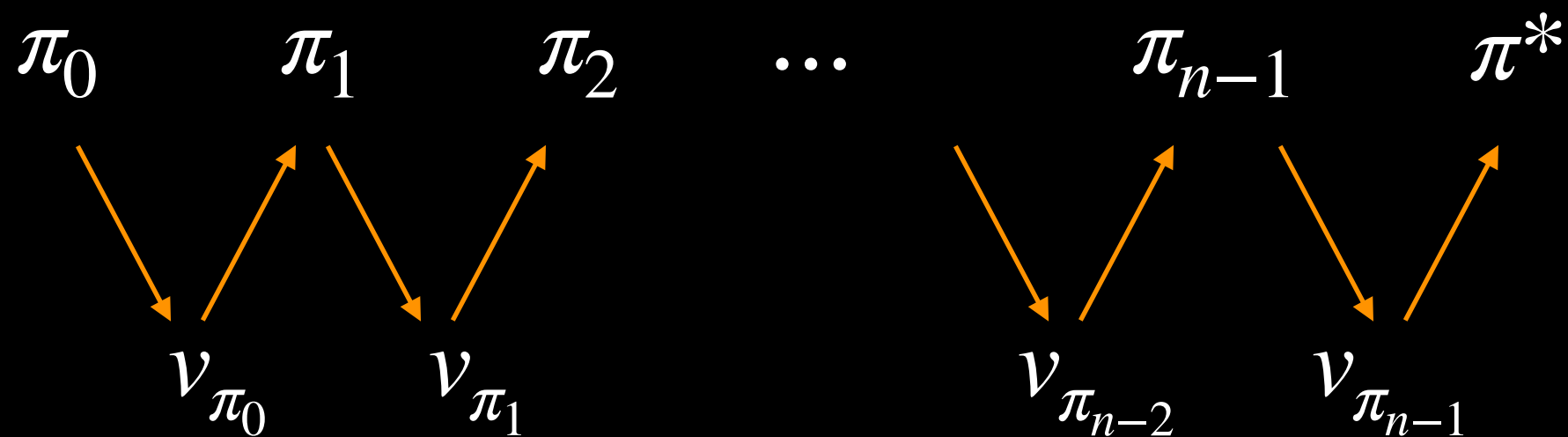
- ▶ Dynamic Programming (DP)
- ▶ Temporal-Difference (TD) Learning
- ▶ Model-based RL
- ▶ Policy Optimization

- ▶ So far we have used value estimates to improve the policy.

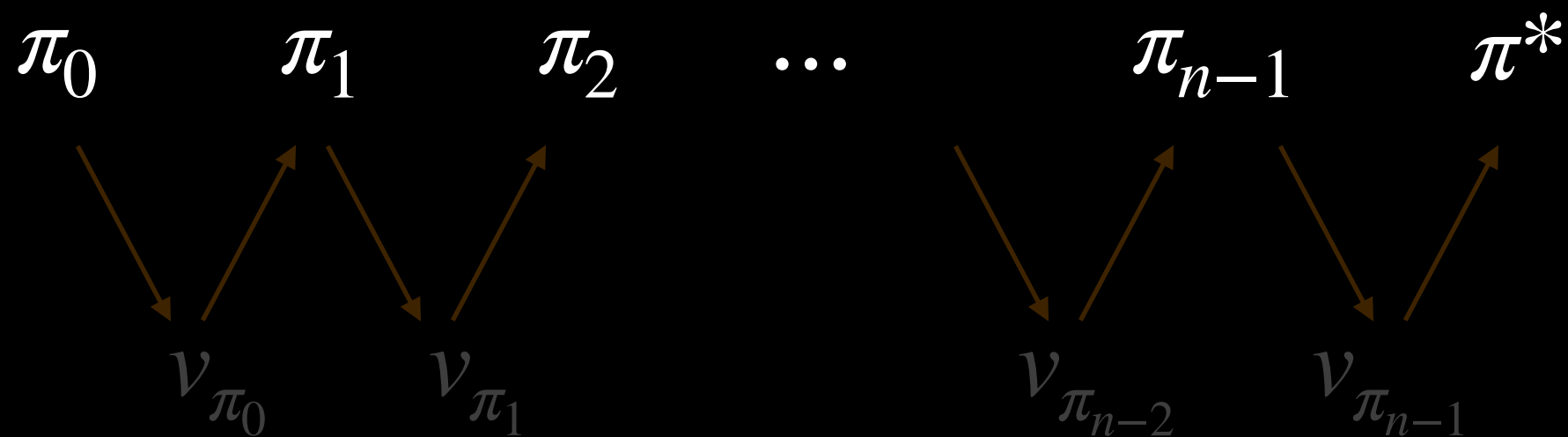
- So far we have used value estimates to improve the policy.



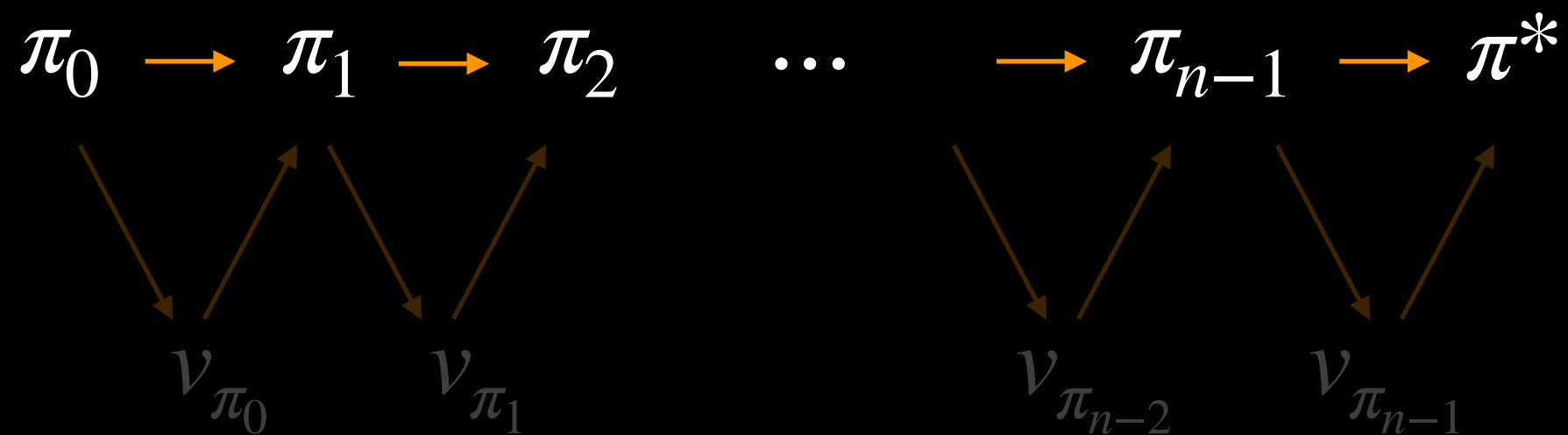
- ▶ So far we have used value estimates to improve the policy.
- ▶ Can we directly improve the policy?



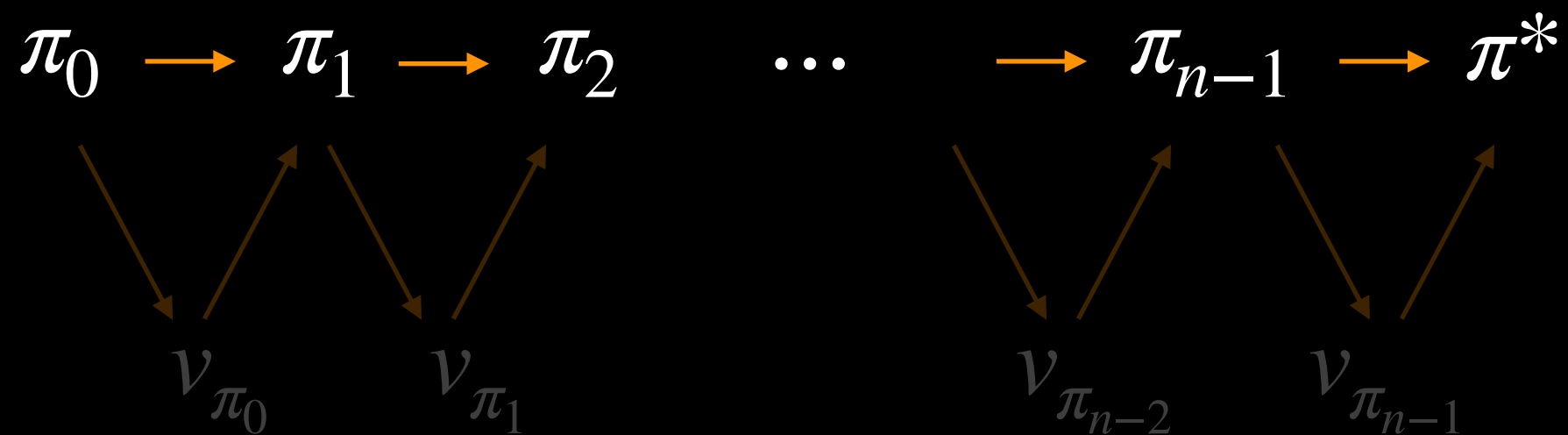
- ▶ So far we have used value estimates to improve the policy.
- ▶ Can we directly improve the policy?



- ▶ So far we have used value estimates to improve the policy.
- ▶ Can we directly improve the policy?



- ▶ So far we have used value estimates to improve the policy.
- ▶ Can we directly improve the policy?



Yes!

POLICY PARAMETRIZATION

POLICY PARAMETRIZATION

- ▶ In value-based methods, we select actions based on the learned value function:

POLICY PARAMETRIZATION

- ▶ In value-based methods, we select actions based on the learned value function:

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

POLICY PARAMETRIZATION

- ▶ In value-based methods, we select actions based on the learned value function:

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

- ▶ In policy-based methods, we need not consult the value function to select actions.

POLICY PARAMETRIZATION

- ▶ In value-based methods, we select actions based on the learned value function:

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

- ▶ In policy-based methods, we need not consult the value function to select actions.
 - ▶ Instead, we use action preferences that are learned.

POLICY PARAMETRIZATION

- ▶ In value-based methods, we select actions based on the learned value function:

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

- ▶ In policy-based methods, we need not consult the value function to select actions.
 - ▶ Instead, we use action preferences that are learned.

$$h(s, a), \forall s, a$$

POLICY PARAMETRIZATION

- ▶ In value-based methods, we select actions based on the learned value function:

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

- ▶ In policy-based methods, we need not consult the value function to select actions.
 - ▶ Instead, we use action preferences that are learned.

$$h(s, a), \forall s, a$$

$$\pi(a | s) = \frac{h(s, a)}{\sum_b h(s, b)}$$

POLICY PARAMETRIZATION

- ▶ In value-based methods, we select actions based on the learned value function:

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

- ▶ In policy-based methods, we need not consult the value function to select actions.
 - ▶ Instead, we use action preferences that are learned.

$$h(s, a), \forall s, a$$

$$\pi(a | s) = \frac{h(s, a)}{\sum_b h(s, b)} \quad \text{or} \quad \frac{e^{h(s, a)}}{\sum_b e^{h(s, b)}}$$

POLICY PARAMETRIZATION

- ▶ In value-based methods, we select actions based on the learned value function:

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

- ▶ In policy-based methods, we need not consult the value function to select actions.
 - ▶ Instead, we use action preferences that are learned.

$$h(s, a), \forall s, a$$

$$\pi(a | s) = \frac{h(s, a)}{\sum_b h(s, b)} \quad \text{or} \quad \frac{e^{h(s, a)}}{\sum_b e^{h(s, b)}}$$

“soft-max”

LEARNING ACTION PREFERENCES

LEARNING ACTION PREFERENCES

- ▶ Consider a simple case: there is just one state.

LEARNING ACTION PREFERENCES

- ▶ Consider a simple case: there is just one state.
- ▶ Objective: maximize amount of reward obtained

LEARNING ACTION PREFERENCES

- ▶ Consider a simple case: there is just one state.
- ▶ Objective: maximize amount of reward obtained
- ▶ Method: stochastic gradient ascent

LEARNING ACTION PREFERENCES

- ▶ Consider a simple case: there is just one state.
- ▶ Objective: maximize amount of reward obtained
- ▶ Method: stochastic gradient ascent

$$\pi(a) = \frac{e^{H(a)}}{\sum_b e^{H(b)}}$$

LEARNING ACTION PREFERENCES

- ▶ Consider a simple case: there is just one state.
- ▶ Objective: maximize amount of reward obtained
- ▶ Method: stochastic gradient ascent

$$\pi(a) = \frac{e^{H(a)}}{\sum_b e^{H(b)}}$$

$$H_{t+1}(A_t) \doteq H_t(A_t) + \alpha(R_t - \bar{R}_t)(1 - \pi(A_t))$$

LEARNING ACTION PREFERENCES

- ▶ Consider a simple case: there is just one state.
- ▶ Objective: maximize amount of reward obtained
- ▶ Method: stochastic gradient ascent

$$\pi(a) = \frac{e^{H(a)}}{\sum_b e^{H(b)}}$$

$$H_{t+1}(A_t) \doteq H_t(A_t) + \alpha(R_t - \bar{R}_t)(1 - \pi(A_t))$$

$$H_{t+1}(a) \doteq H_t(a) + \alpha(R_t - \bar{R}_t)(0 - \pi(a)) \quad \forall a \neq A_t$$

LEARNING ACTION PREFERENCES

- ▶ Consider a simple case: there is just one state.
- ▶ Objective: maximize amount of reward obtained
- ▶ Method: stochastic gradient ascent

$$\pi(a) = \frac{e^{H(a)}}{\sum_b e^{H(b)}}$$

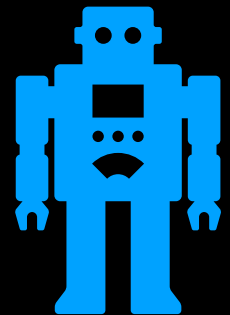
“Gradient-Bandit Algorithm”

$$H_{t+1}(A_t) \doteq H_t(A_t) + \alpha(R_t - \bar{R}_t)(1 - \pi(A_t))$$

$$H_{t+1}(a) \doteq H_t(a) + \alpha(R_t - \bar{R}_t)(0 - \pi(a)) \quad \forall a \neq A_t$$

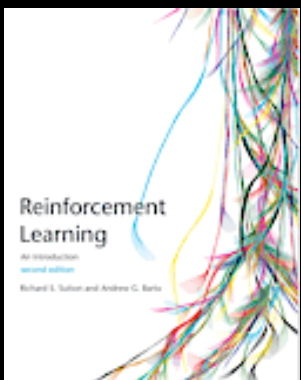
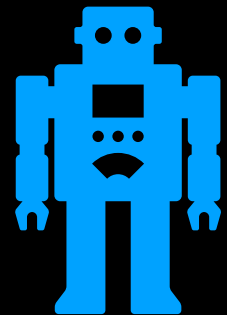
OUTLINE

- ▶ Dynamic Programming (DP)
- ▶ Temporal-Difference (TD) Learning
- ▶ Model-based RL
- ▶ Policy Optimization



OUTLINE

- ▶ Dynamic Programming (DP)
- ▶ Temporal-Difference (TD) Learning
- ▶ Model-based RL
- ▶ Policy Optimization



Further reading:

Sutton & Barto, 2018, *Reinforcement Learning: An Introduction*, 2nd Edition

<http://incompleteideas.net/book/the-book.html>

THANK YOU

Questions?

 abhisheknaik96.github.io

 abhishek.naik@ualberta.ca

 [anaik96](https://twitter.com/anaik96)