

Computational Sensorimotor Learning (Spring'21)

Pulkit Agrawal

Lecture 3
Feb 23 2021

Consider landing pages for your website

 a_1 a_2 a_3

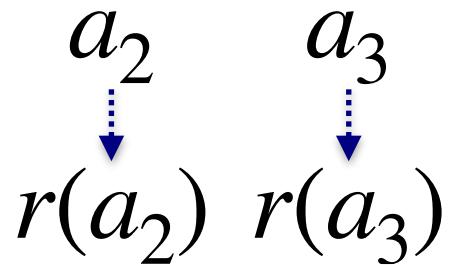

$$a_2 \quad a_3$$

$r(a_2) \quad r(a_3)$


$$a_1 \quad a_2 \quad a_3$$



(male, 30s, computer-savvy)



(female, 20s, computer-savvy)

How to use these “features” in decision making?

Contextual Bandits

LinUCB

a_1

a_2

a_3



$$a_2 \quad a_3$$
$$\downarrow \quad \downarrow$$
$$r(a_2) \ r(a_3)$$



(male, 30s, computer-savvy)

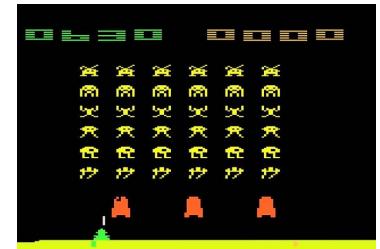
(female, 20s, computer-savvy)

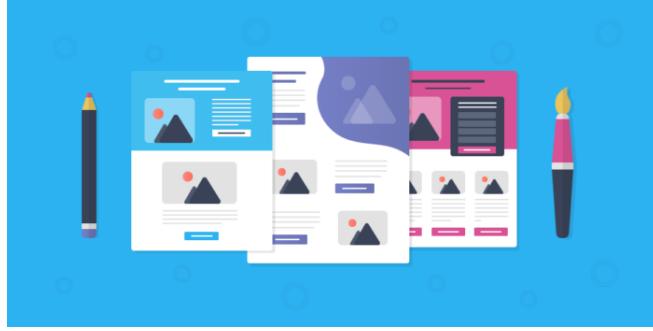
How to use these “features” in decision making?

Contextual Bandits

BUT, Actions don't change the future

Reinforcement Learning

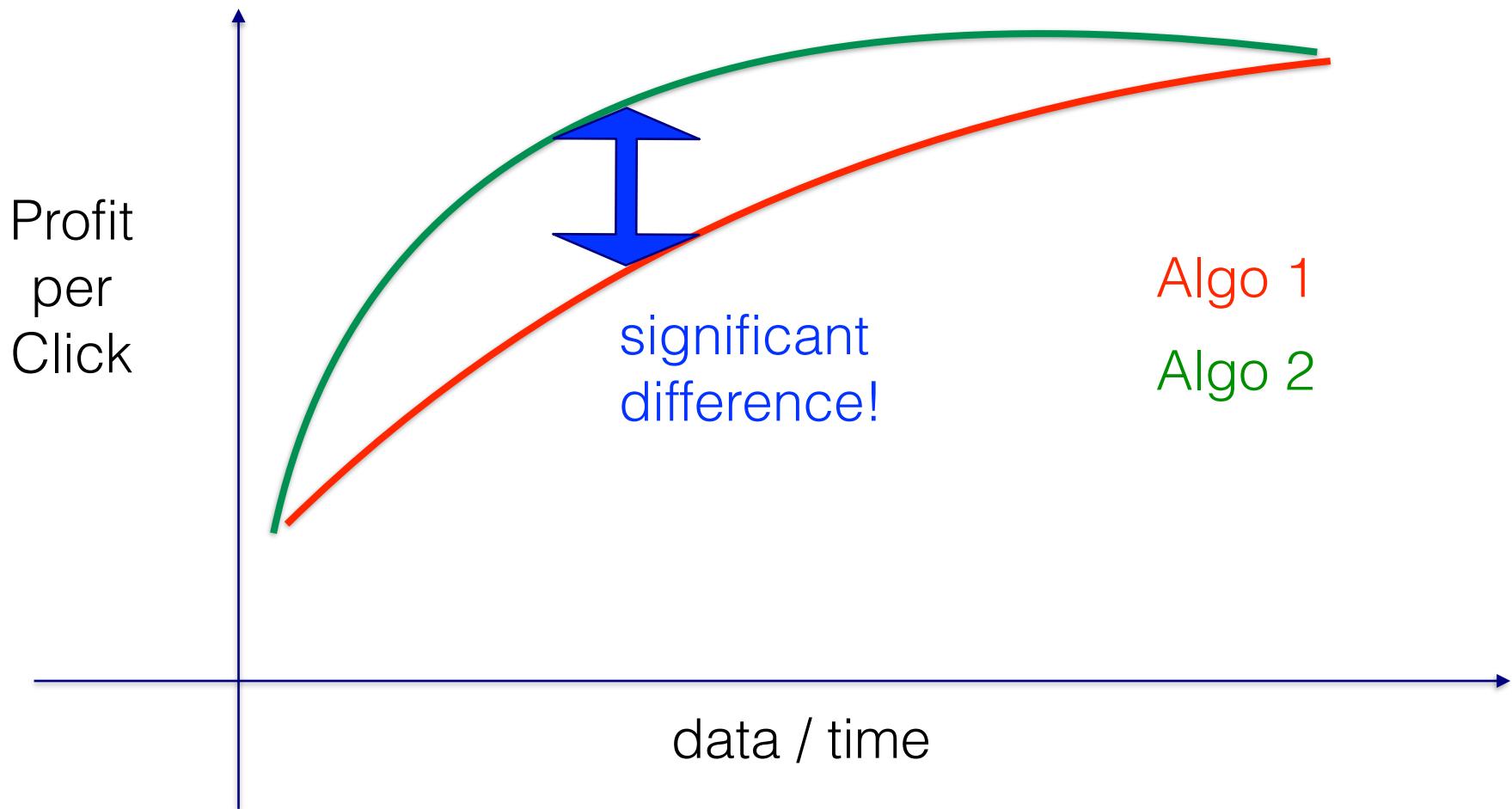




Online
Decisions!

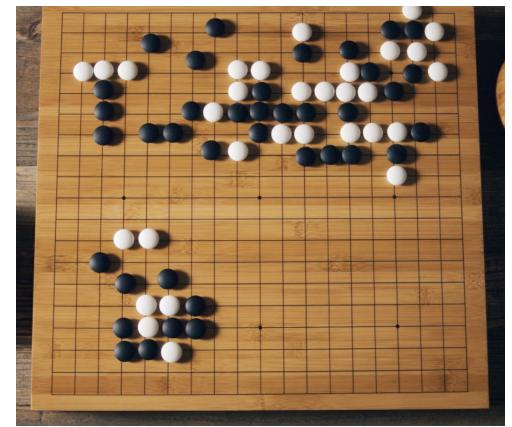
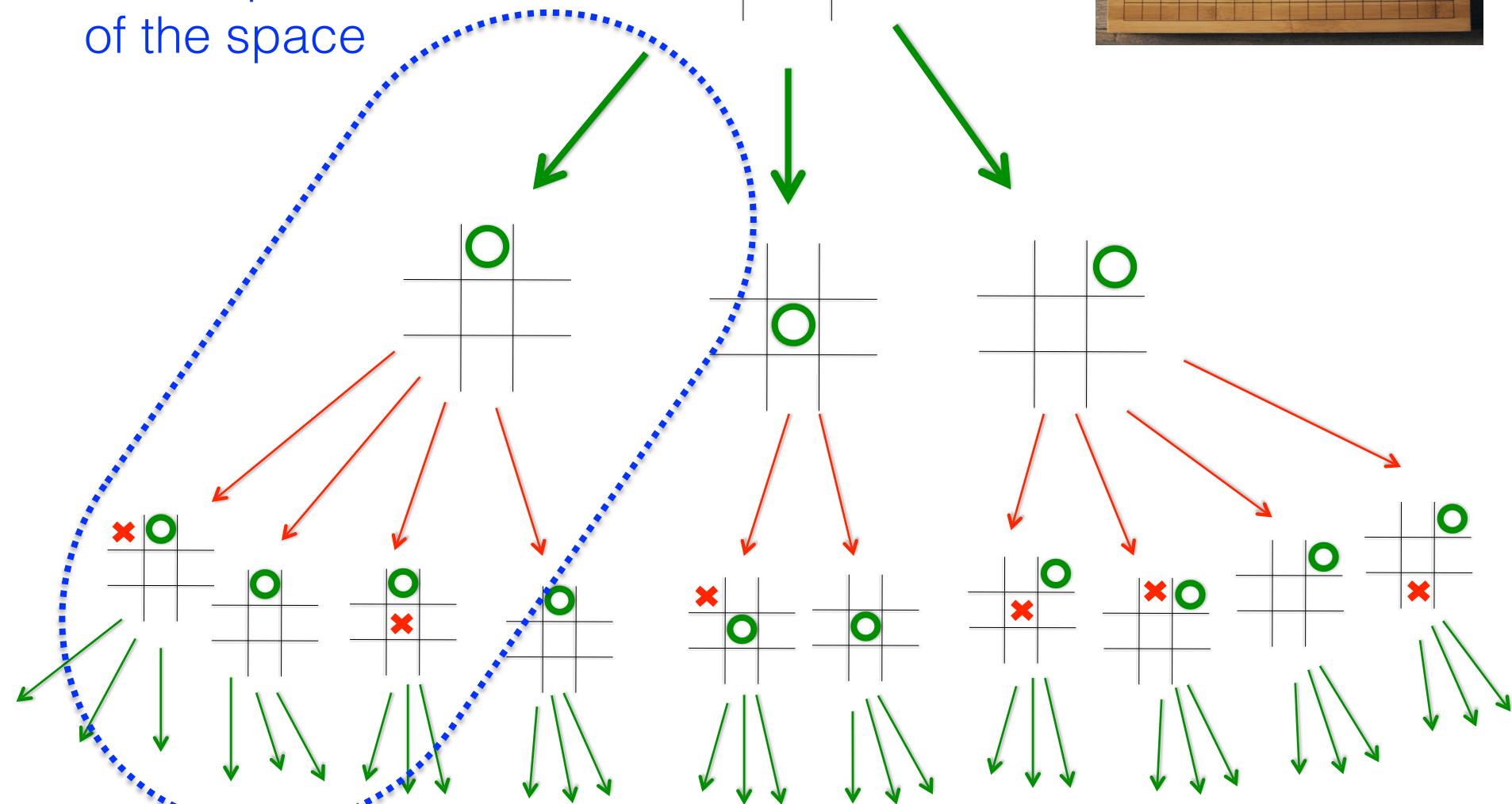
When to use RL?

(i.e., can't wait to collect data first)



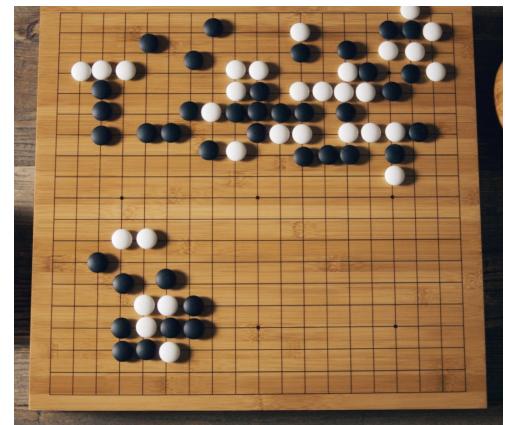
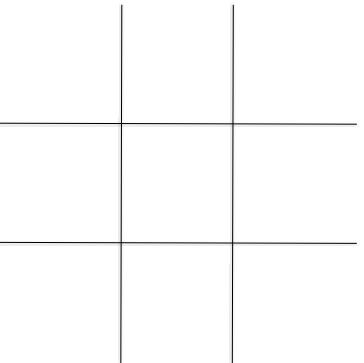
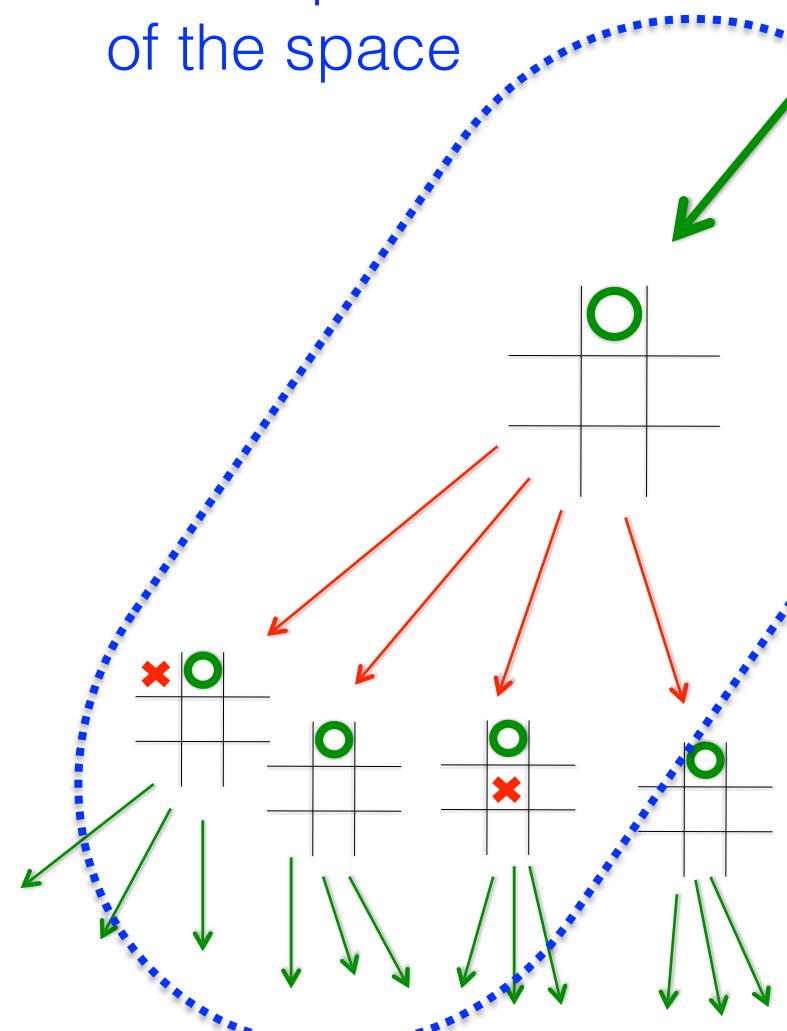
When to use RL?

can only explore
small part
of the space



When to use RL?

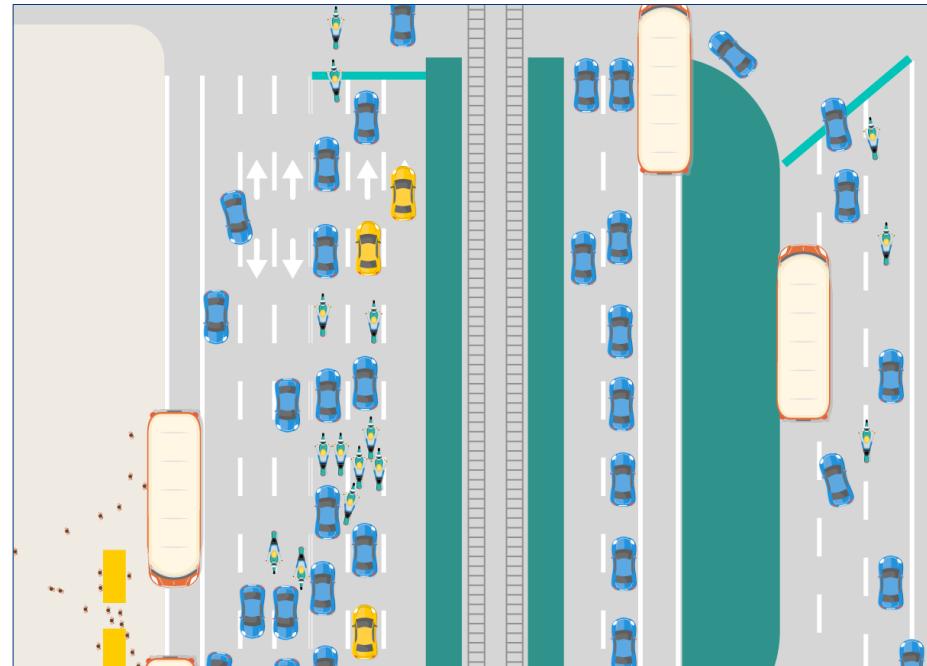
can only explore
small part
of the space



Can we “intelligently”
explore the space to create
new insights?

Knowledge Synthesis!

Example of Knowledge Synthesis: Urban Planning



We have simulators

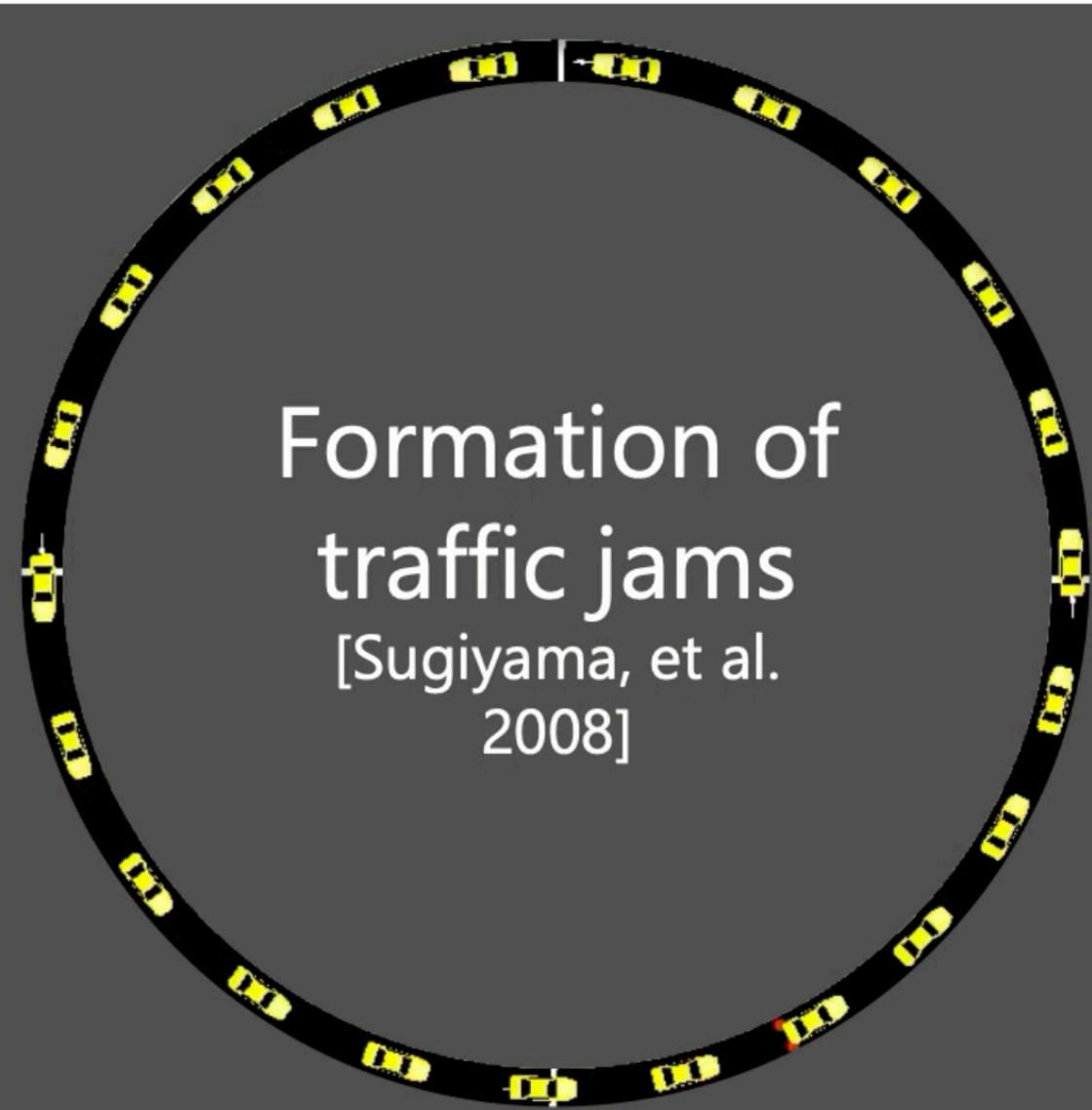
But how to use them for the desired purpose?

Traffic Jam Problem



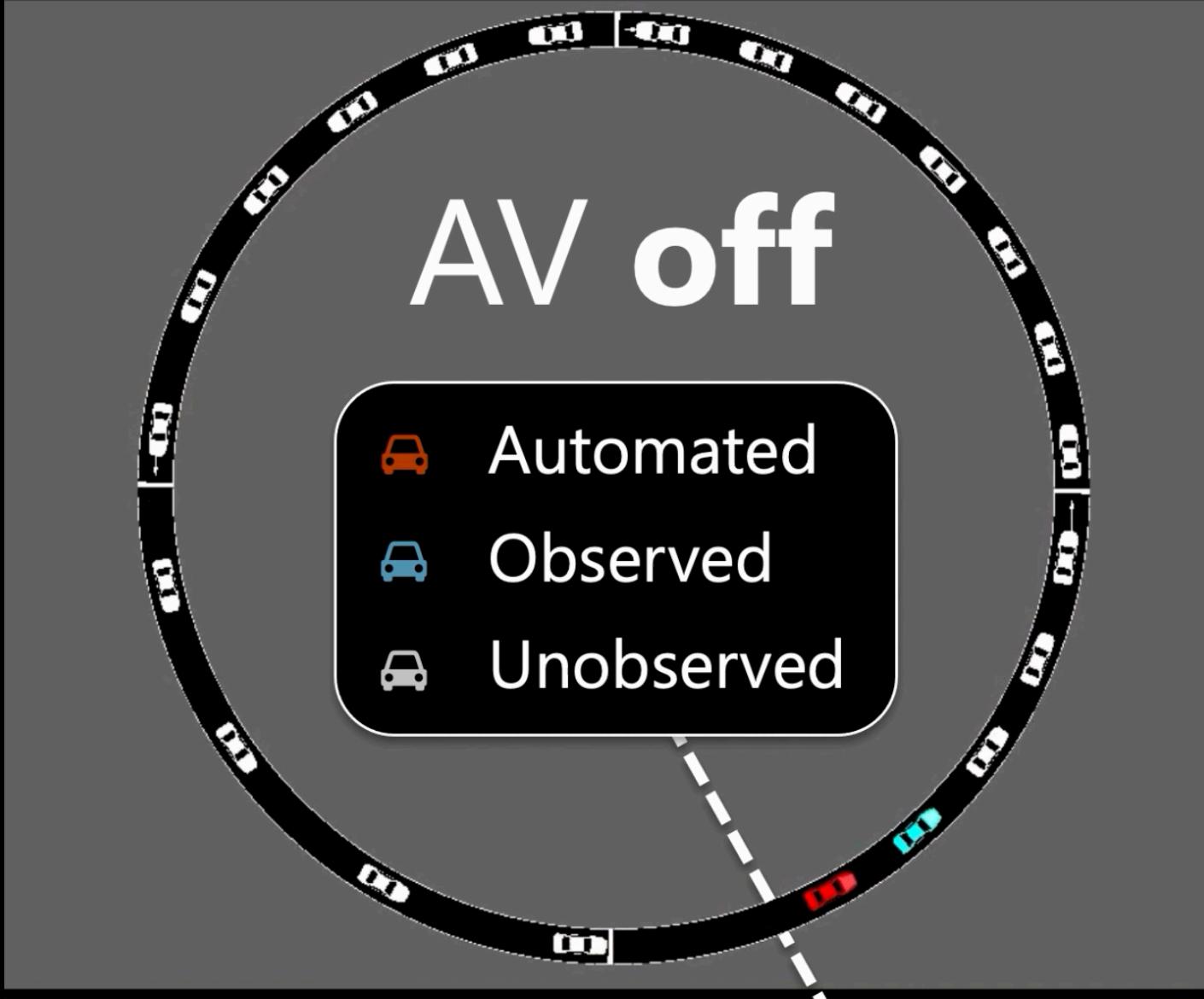
The Mathematical Society of Traffic Flow

Simulated Illustration



The image shows a circular simulation visualization of a ring road. The road is black with yellow dashed lines and features small yellow car icons placed at regular intervals along its circumference. In the center of the circle, white text reads "Formation of traffic jams [Sugiyama, et al. 2008]".

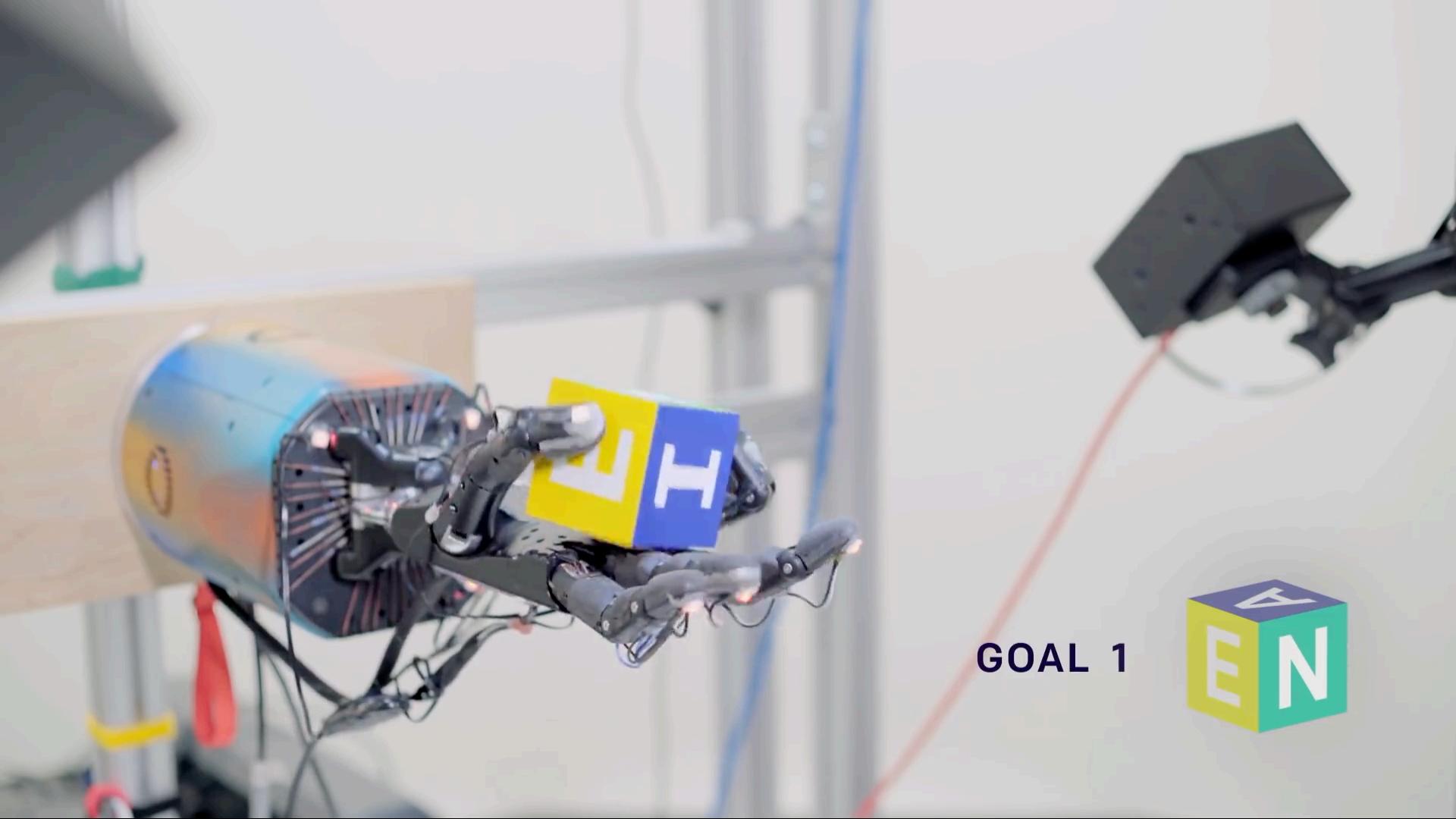
Formation of
traffic jams
[Sugiyama, et al.
2008]



AV off

- Automated
- Observed
- Unobserved

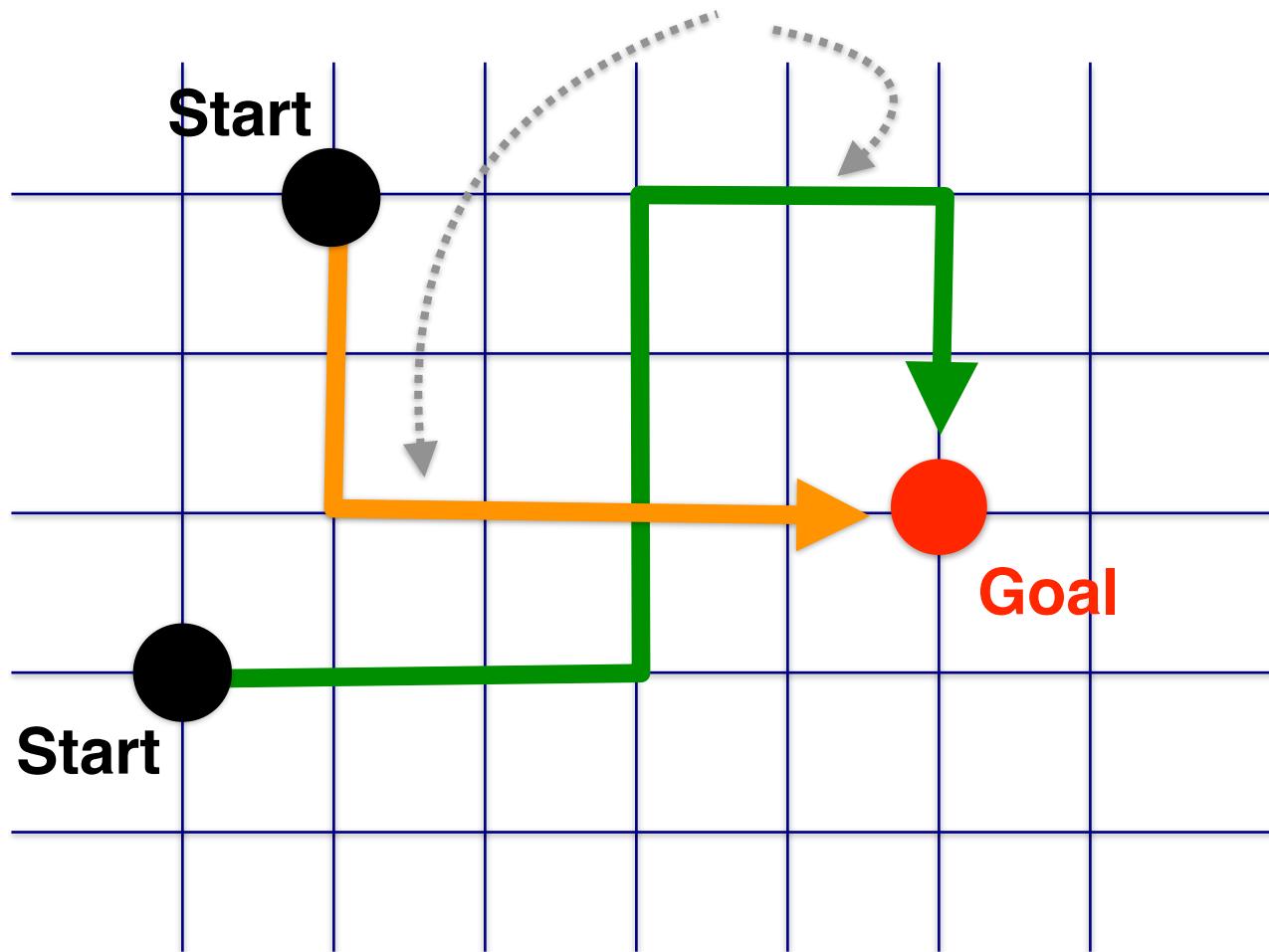
Learning Dexterity



Known Simulator, not differentiable, 24 degrees of freedom!

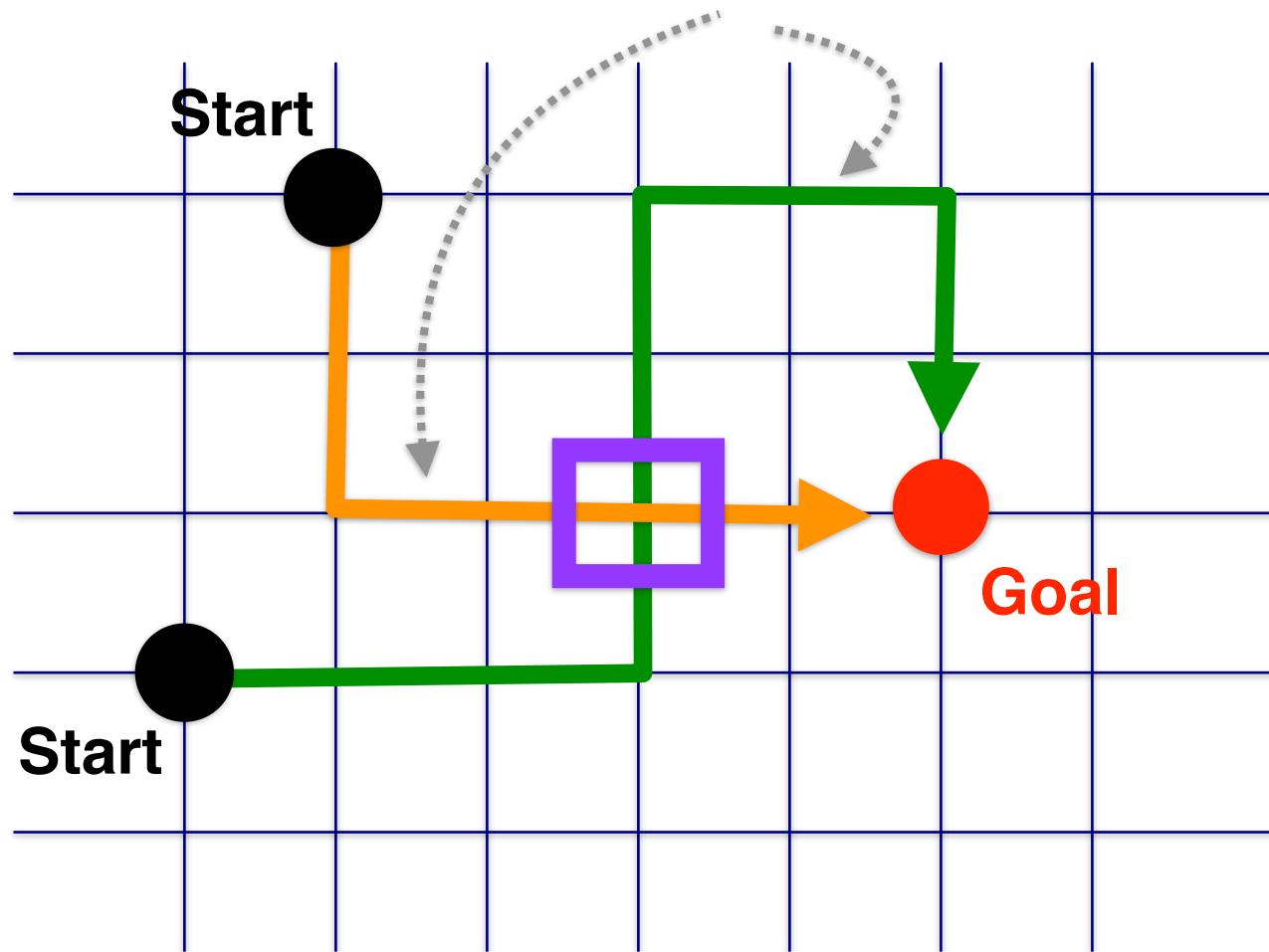
When to use RL?

Existing Data / Solution



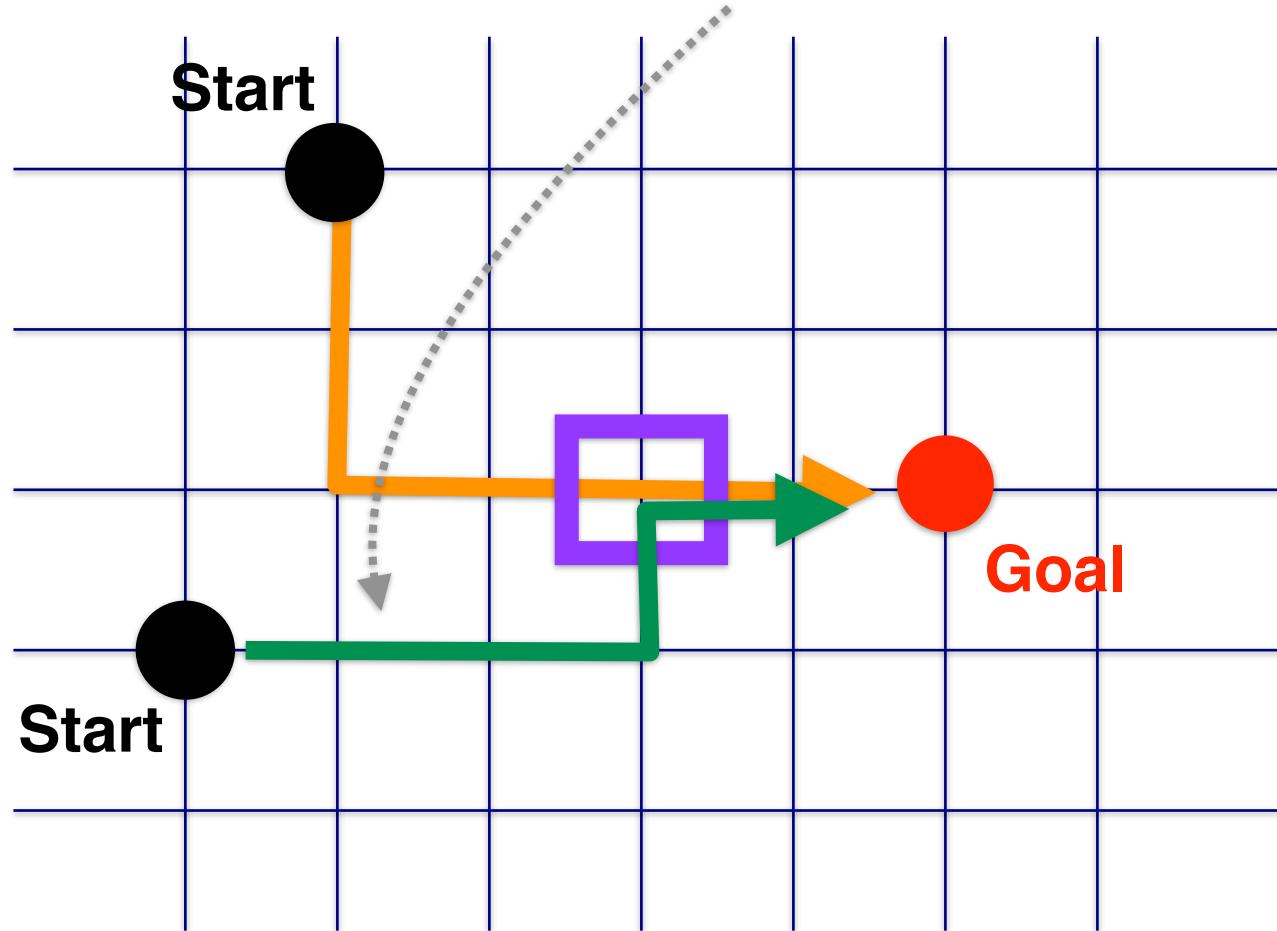
When to use RL?

Existing Data / Solution



When to use RL?

Improve the solution!

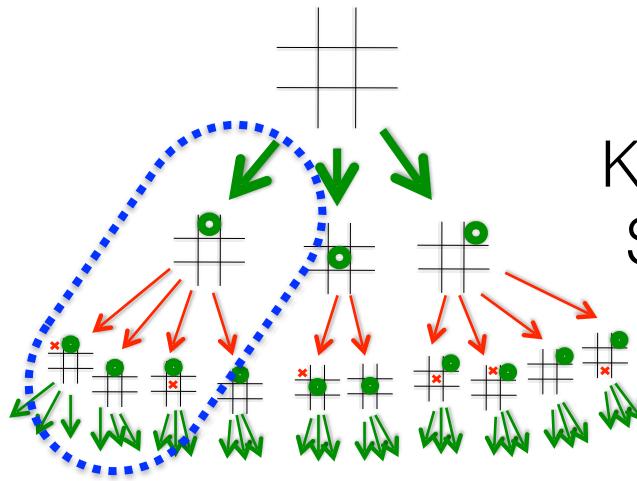


RL can be thought of as finding the “shortest path”

When to use RL?

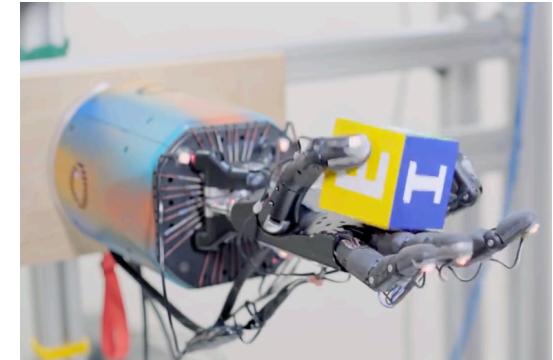


online decisions

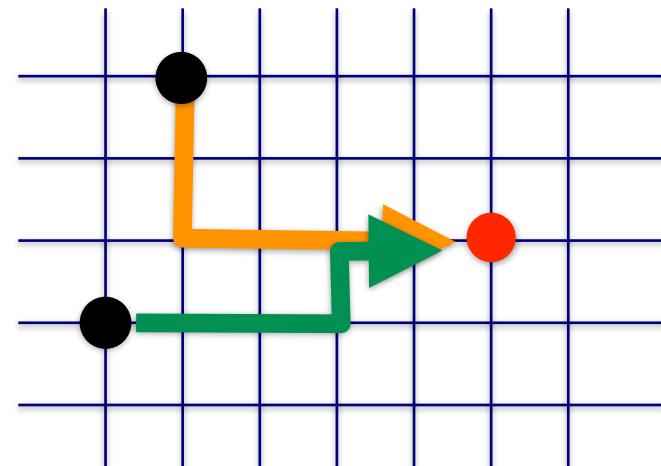


Knowledge
Synthesis

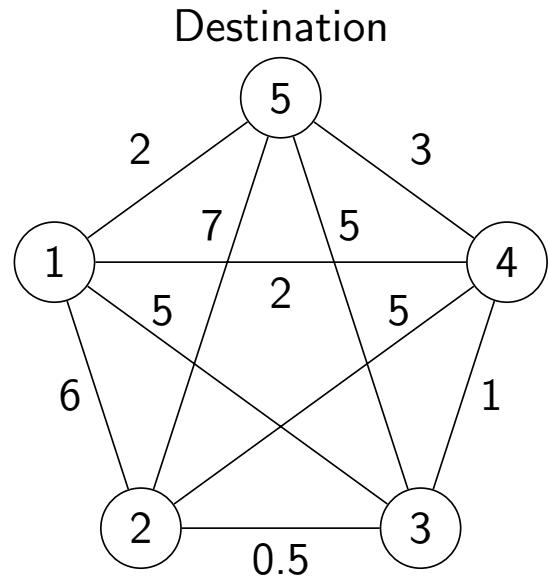
Improve existing solutions



Control non-linear
systems

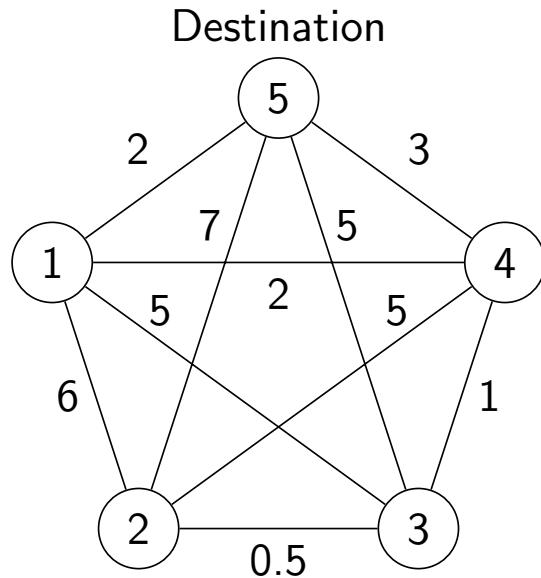


Understanding RL via Shortest Path Problem



Destination is node 5.

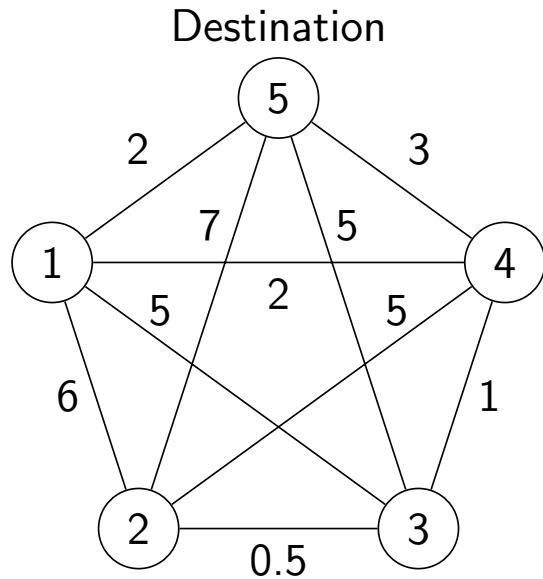
Understanding RL via Shortest Path Problem



Sequential decision problem

Destination is node 5.

Understanding RL via Shortest Path Problem

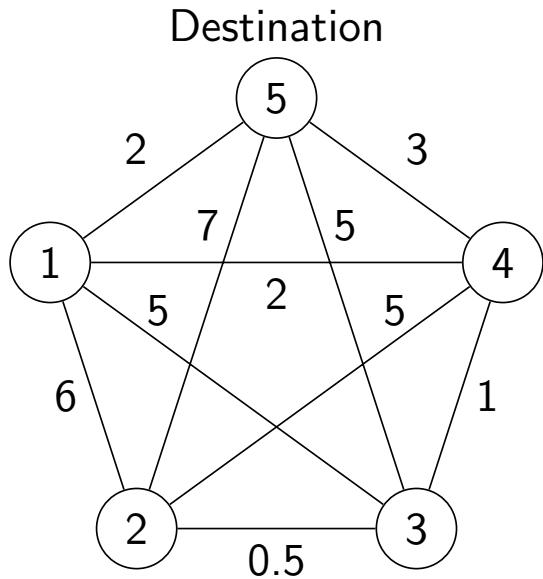


Sequential decision problem

- Start state s_0 : city 2

Destination is node 5.

Understanding RL via Shortest Path Problem

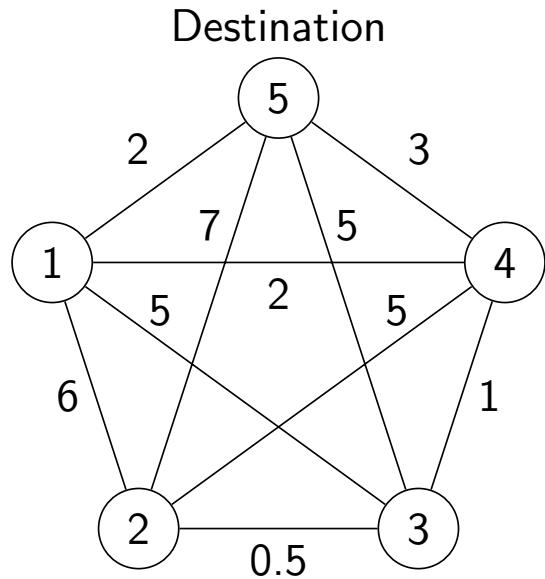


Sequential decision problem

- Start state s_0 : city 2
- Action a_0 : take link between city 2 and city 3

Destination is node 5.

Understanding RL via Shortest Path Problem

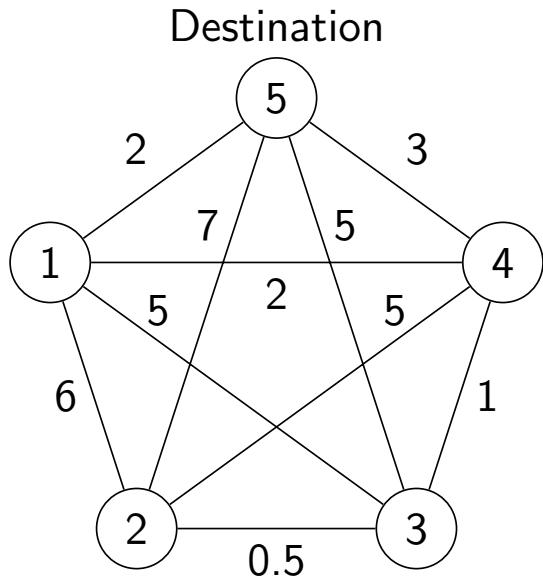


Sequential decision problem

- Start state s_0 : city 2
- Action a_0 : take link between city 2 and city 3
- State s_1 : city 3

Destination is node 5.

Understanding RL via Shortest Path Problem

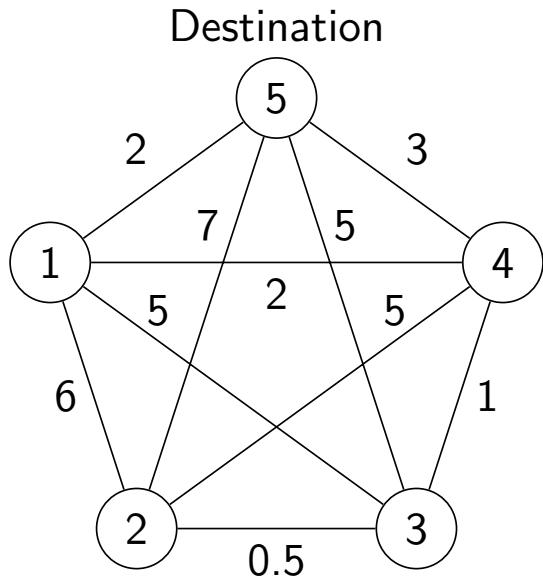


Sequential decision problem

- Start state s_0 : city 2
- Action a_0 : take link between city 2 and city 3
- State s_1 : city 3
- Action a_1 : take link between city 3 and city 5

Destination is node 5.

Understanding RL via Shortest Path Problem

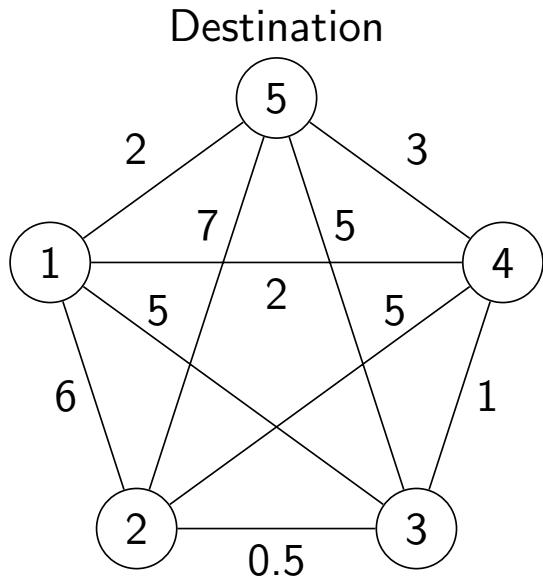


Sequential decision problem

- Start state s_0 : city 2
- Action a_0 : take link between city 2 and city 3
- State s_1 : city 3
- Action a_1 : take link between city 3 and city 5
- State s_2 : city 5

Destination is node 5.

Understanding RL via Shortest Path Problem

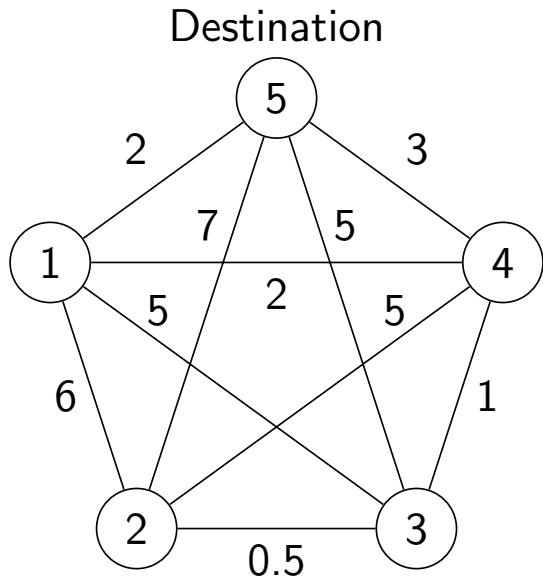


Sequential decision problem

- Start state s_0 : city 2
- Action a_0 : take link between city 2 and city 3
- State s_1 : city 3
- Action a_1 : take link between city 3 and city 5
- State s_2 : city 5
- ...

Destination is node 5.

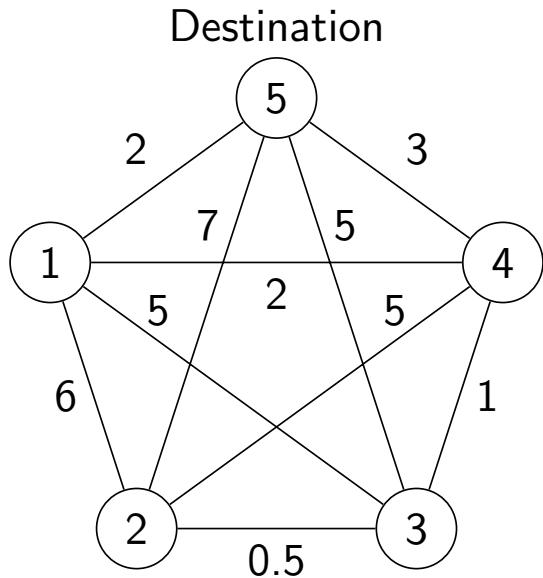
Solving the Shortest Path Problem



- **Naive approach:** enumerate all possibilities.

Destination is node 5.

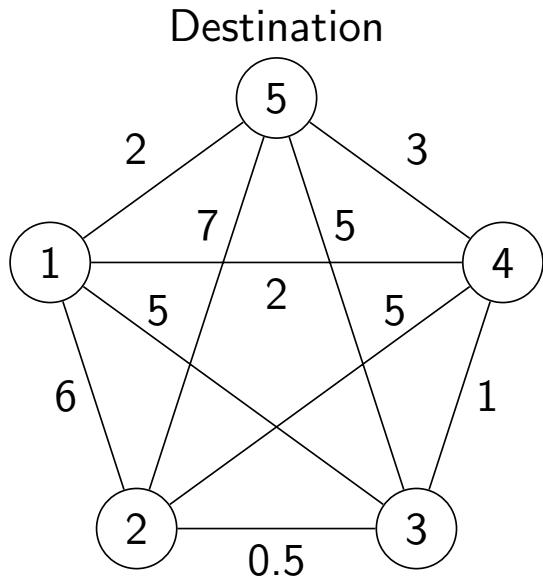
Solving the Shortest Path Problem



- **Naive approach:** enumerate all possibilities.
 - From a starting city s_0 , choose any remaining city ($N - 1$ choices). Choose any next remaining city ($N - 2$ choices). ... Until there is only 1 option remaining.
 - Add up the edge costs.
 - Select the best sequence (lowest total cost).
 - $\mathcal{O}(N!)$. ☹

Destination is node 5.

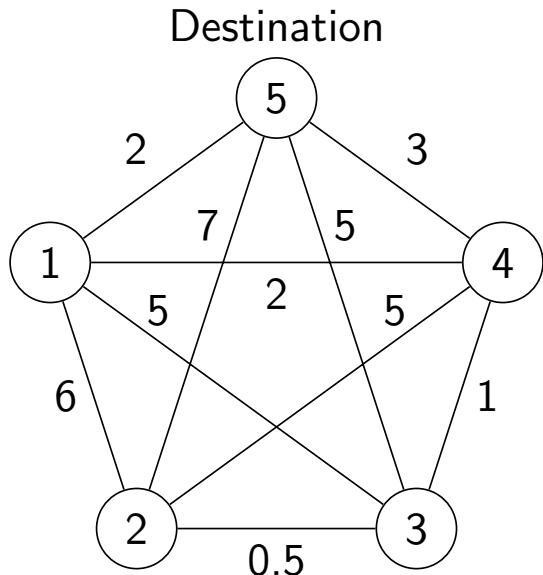
Solving the Shortest Path Problem



- **Issue:** repeated calculations of subsequences.
- Dynamic programming

Destination is node 5.

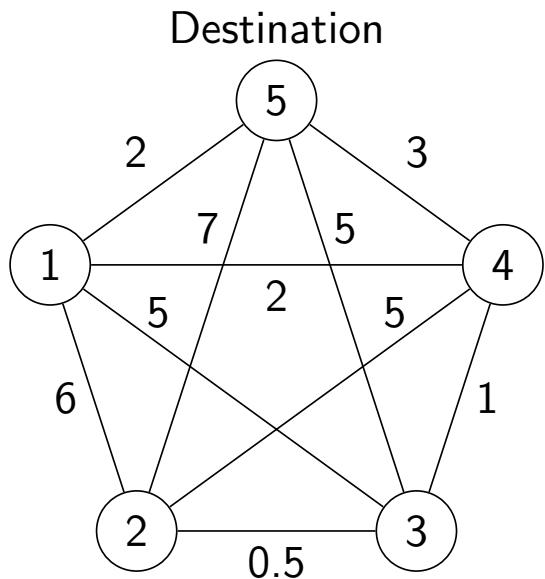
Solving the Shortest Path Problem



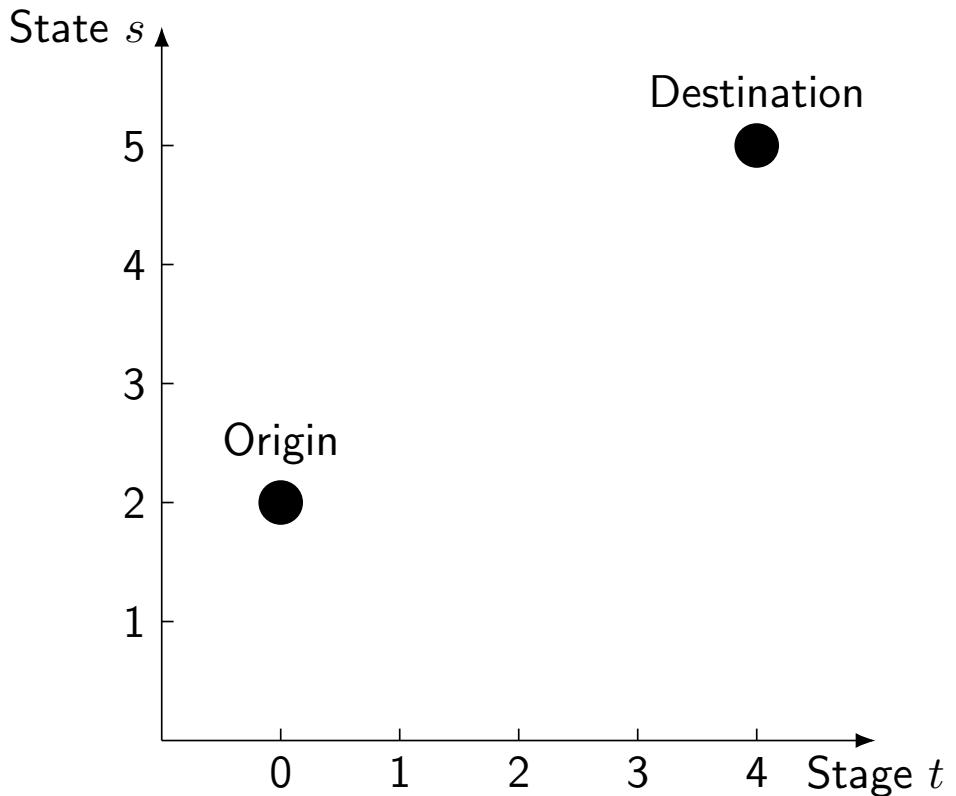
Destination is node 5.

- **Issue:** repeated calculations of subsequences.
- Dynamic programming
 - Divide-and-conquer, or [the principle of optimality](#).
 - Overall problem would be much easier to solve if a part of the problem were already solved.
 - Break a problem down into subproblems.

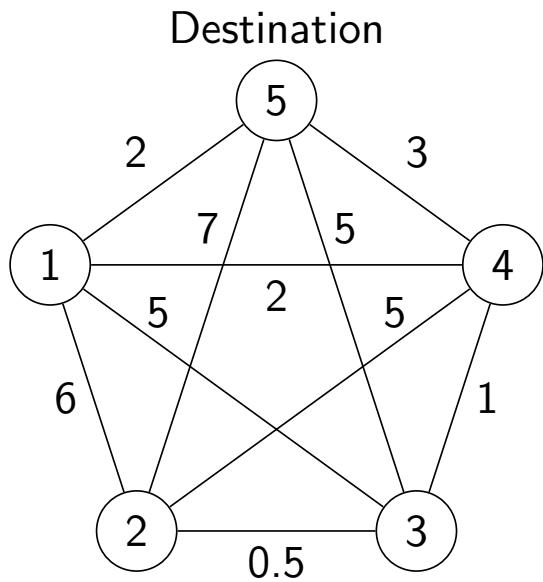
Solving the Shortest Path Problem



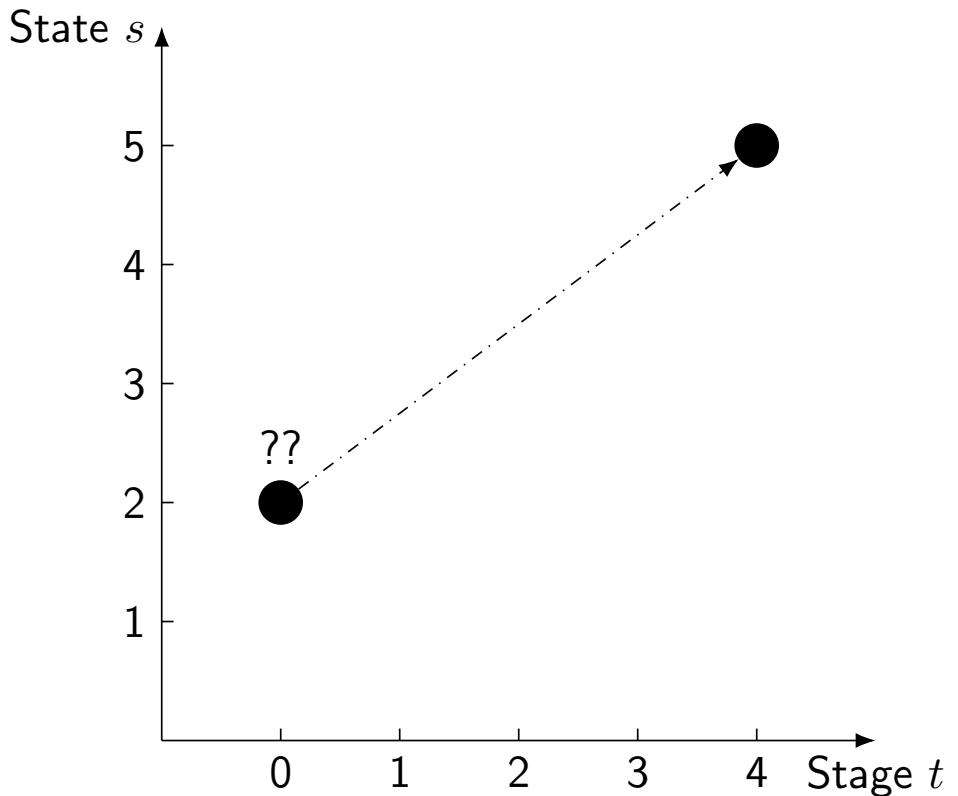
Destination is node 5.



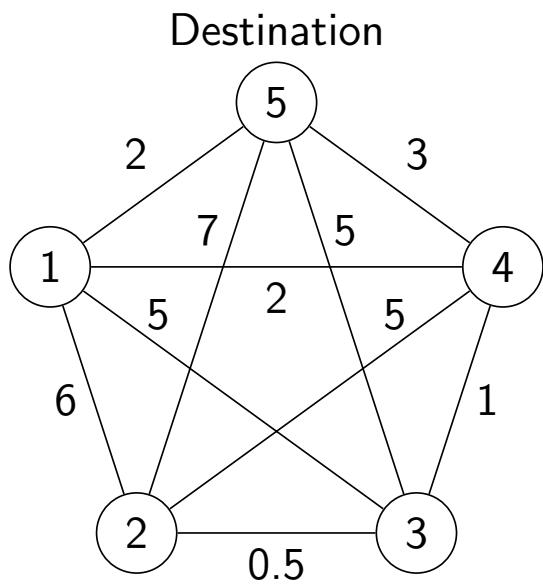
Solving the Shortest Path Problem



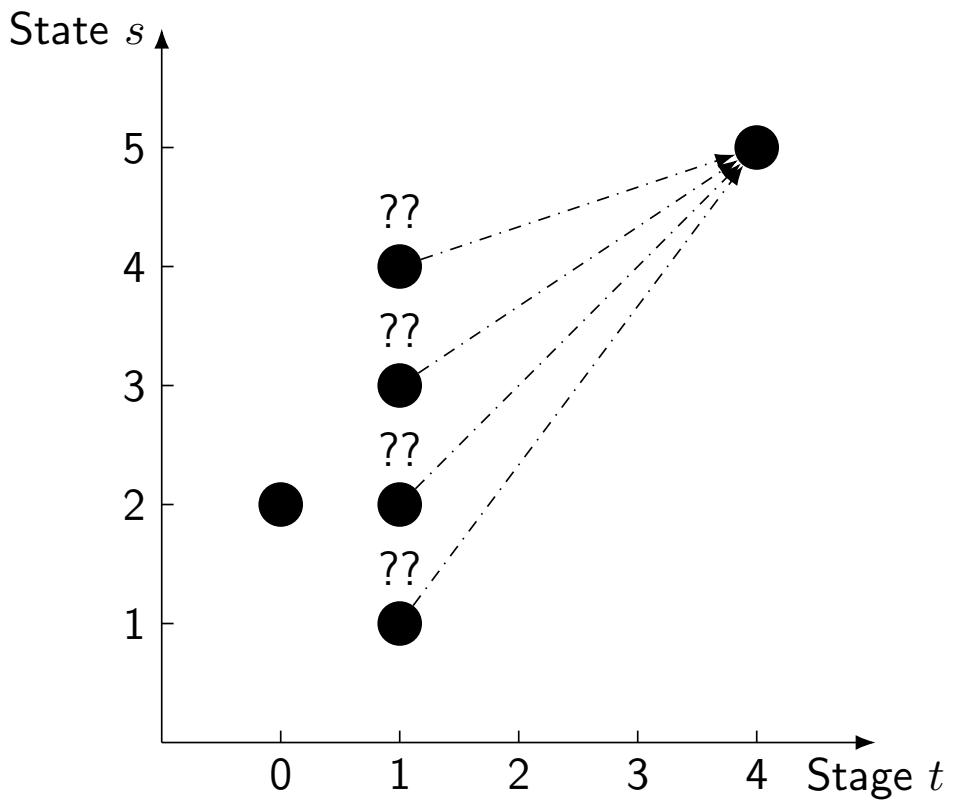
Destination is node 5.



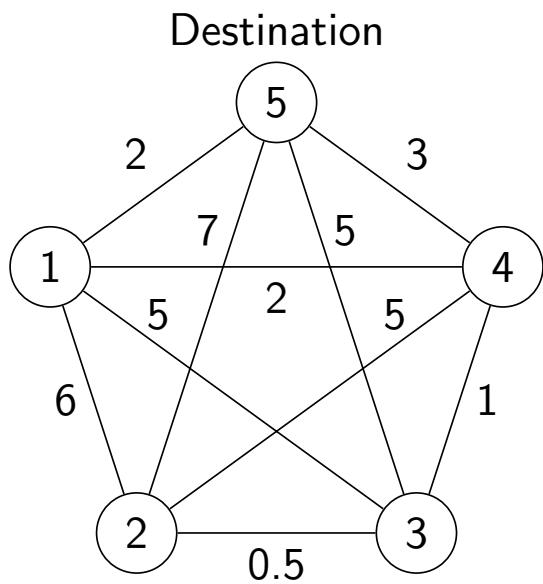
Solving the Shortest Path Problem



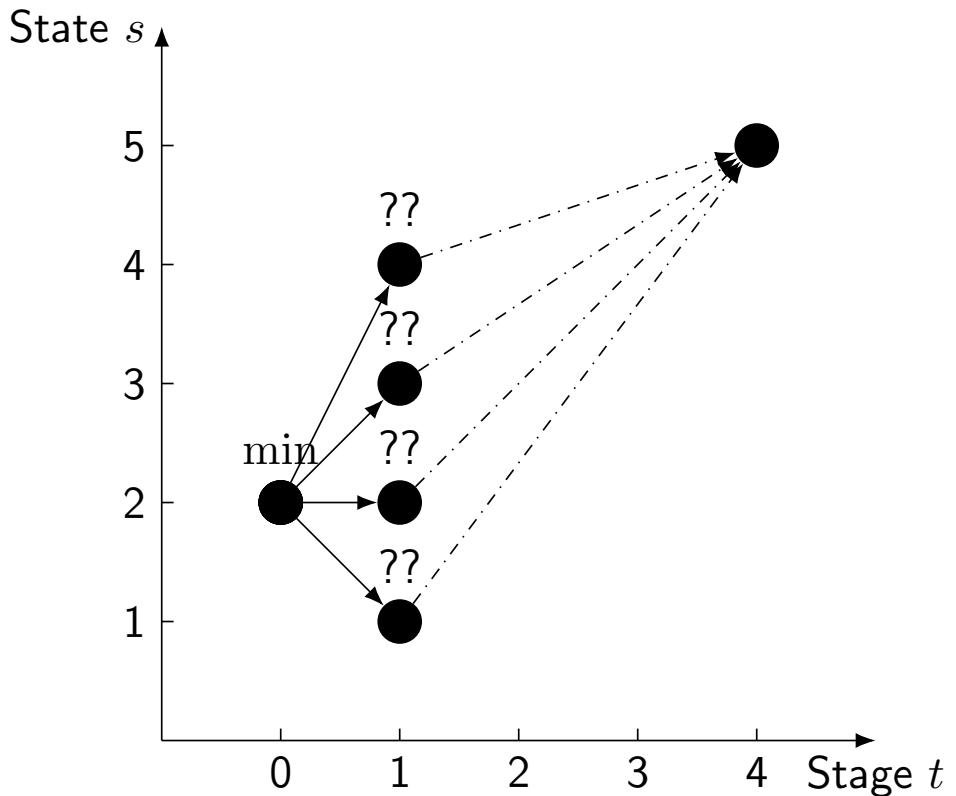
Destination is node 5.



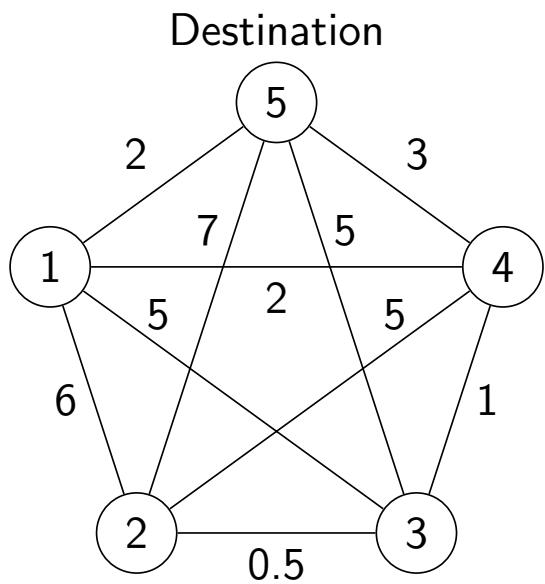
Solving the Shortest Path Problem



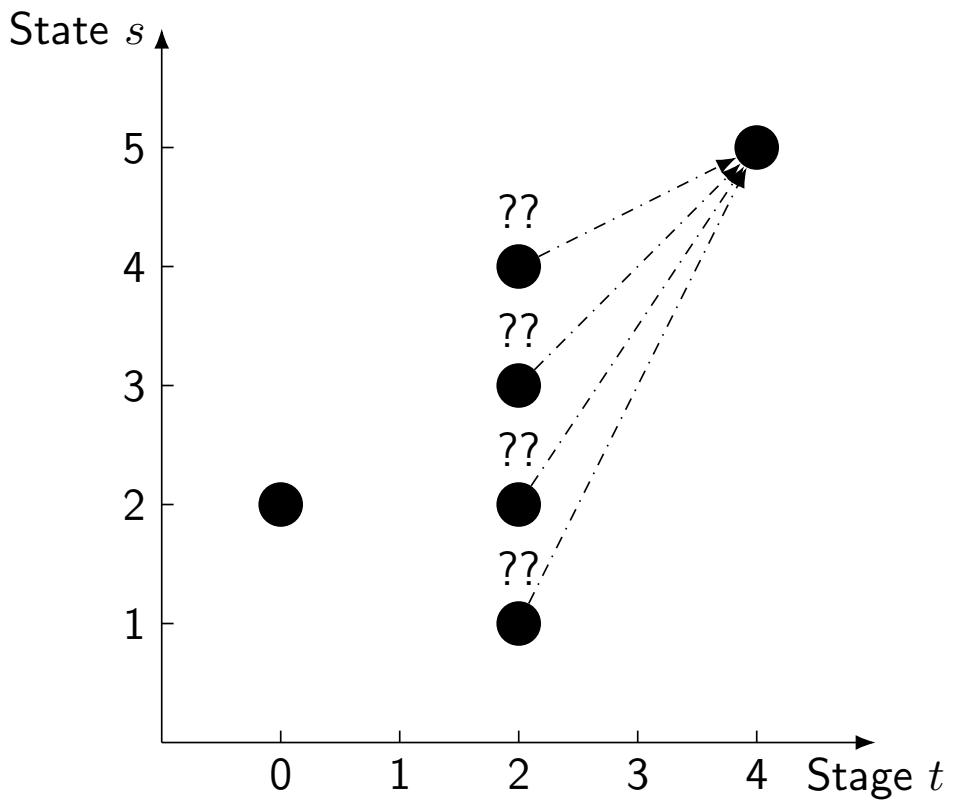
Destination is node 5.



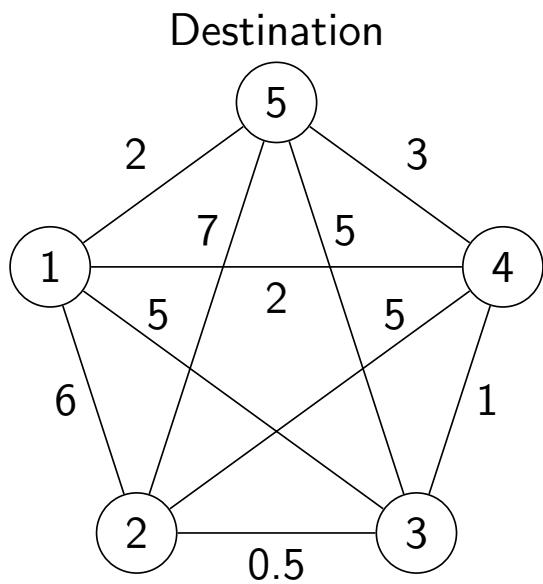
Solving the Shortest Path Problem



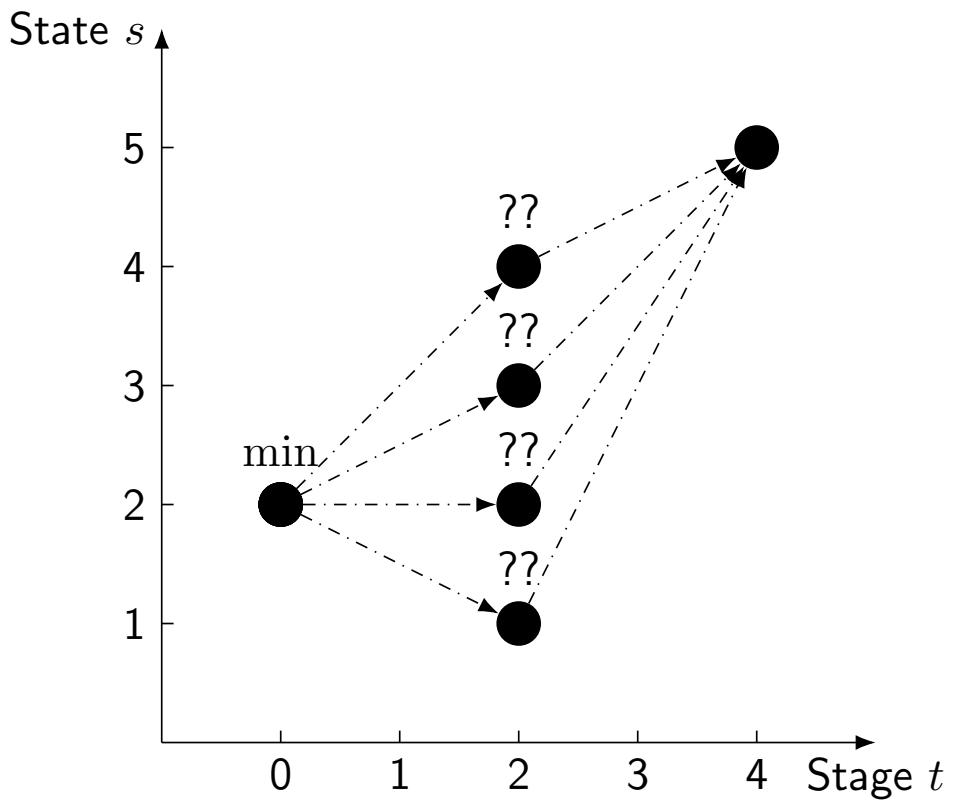
Destination is node 5.



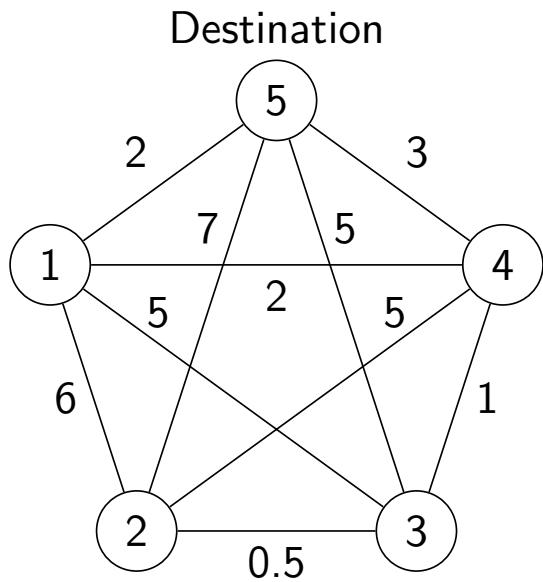
Solving the Shortest Path Problem



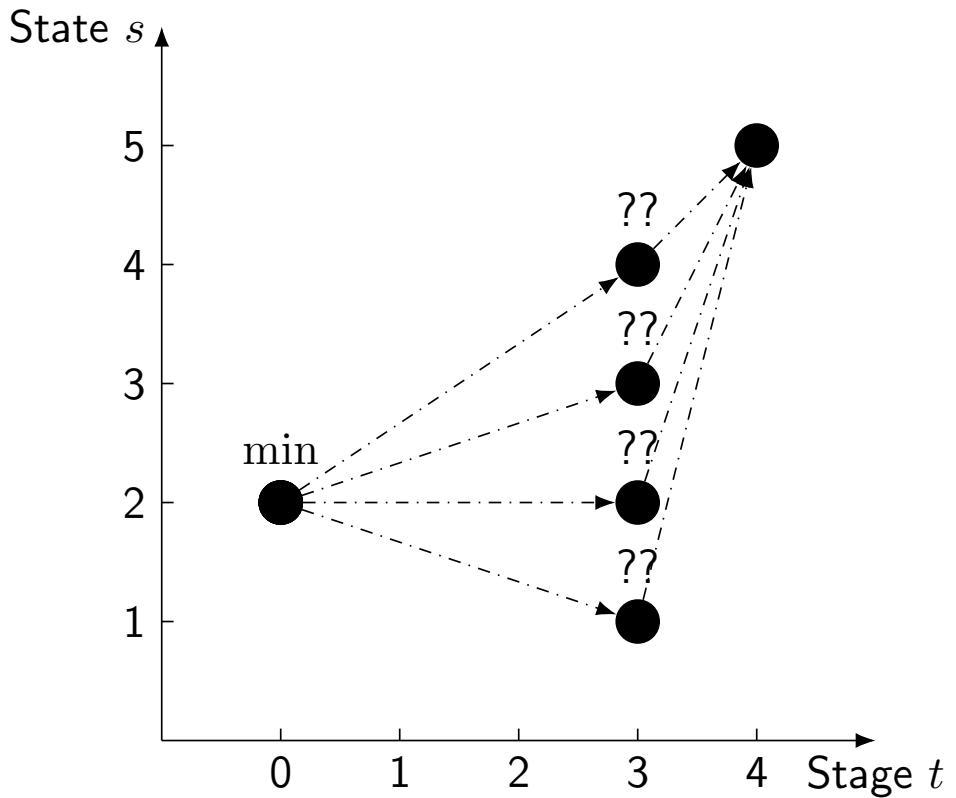
Destination is node 5.



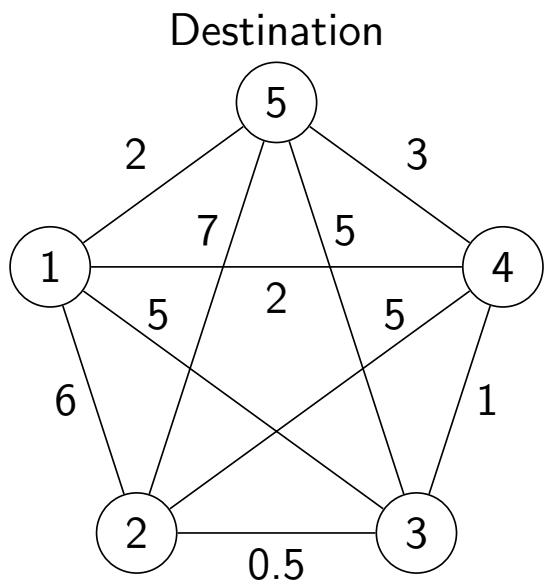
Solving the Shortest Path Problem



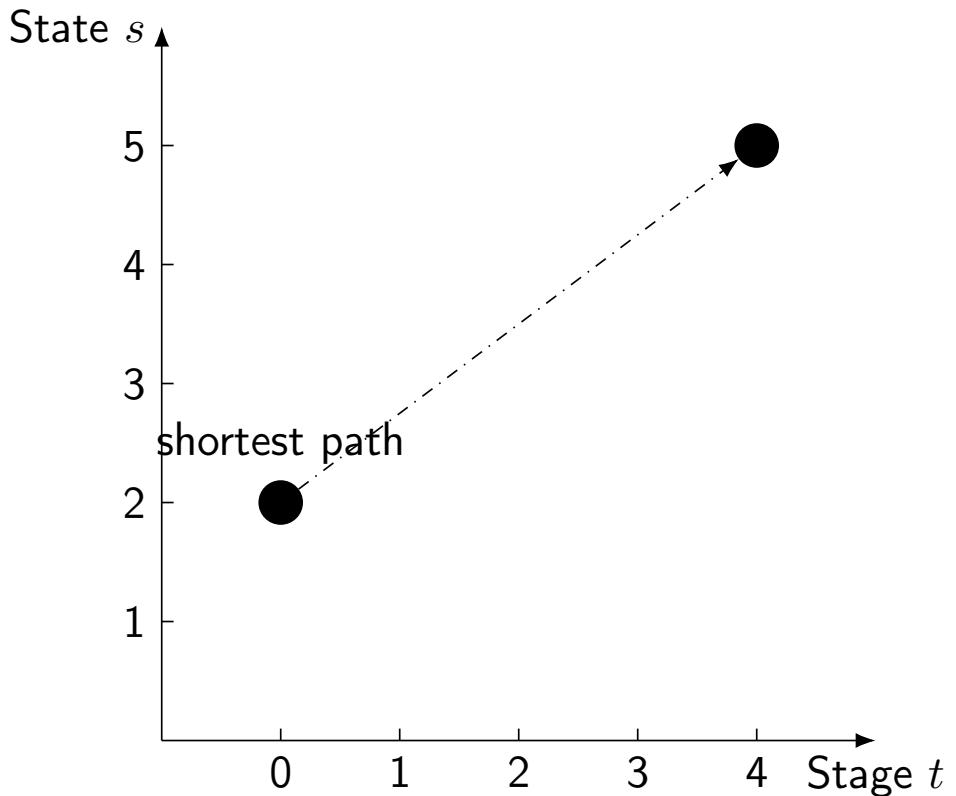
Destination is node 5.



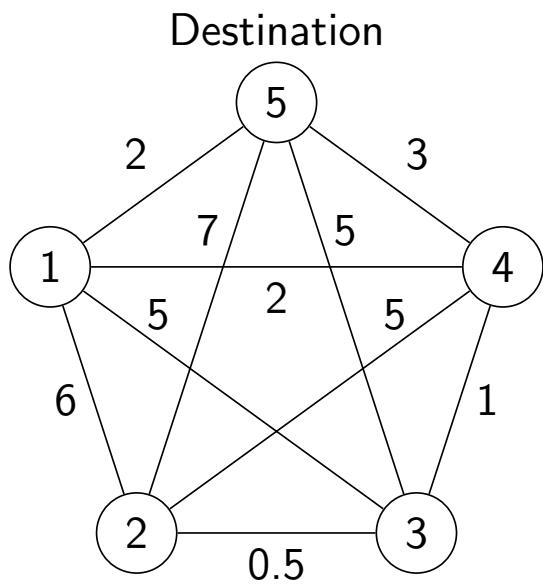
Solving the Shortest Path Problem



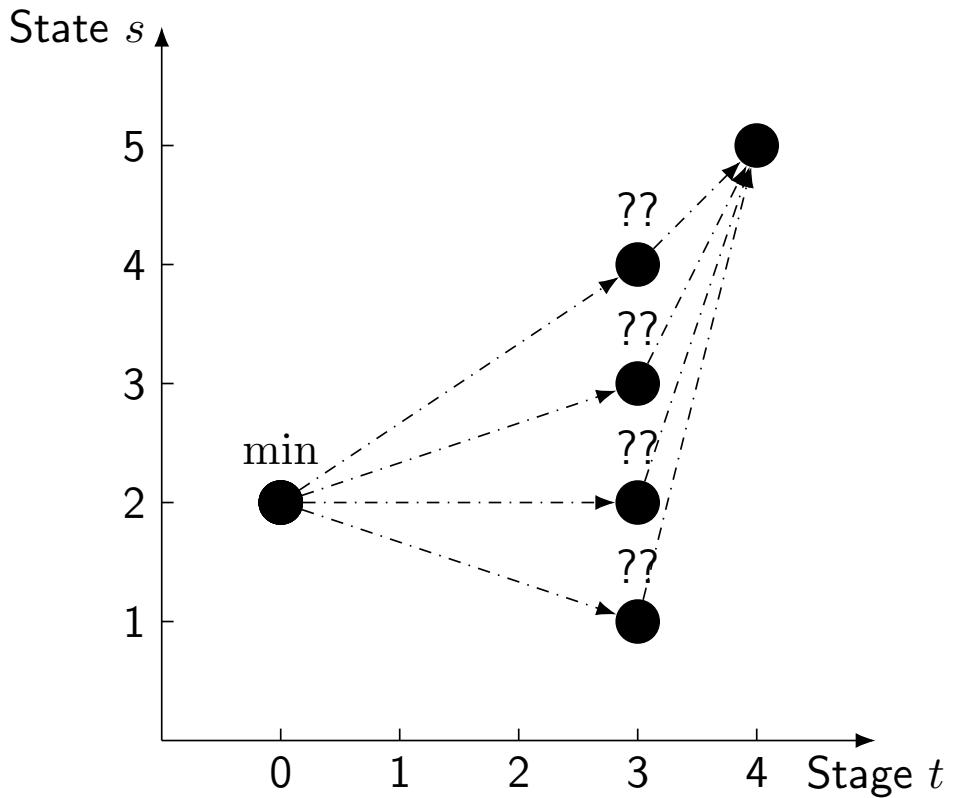
Destination is node 5.



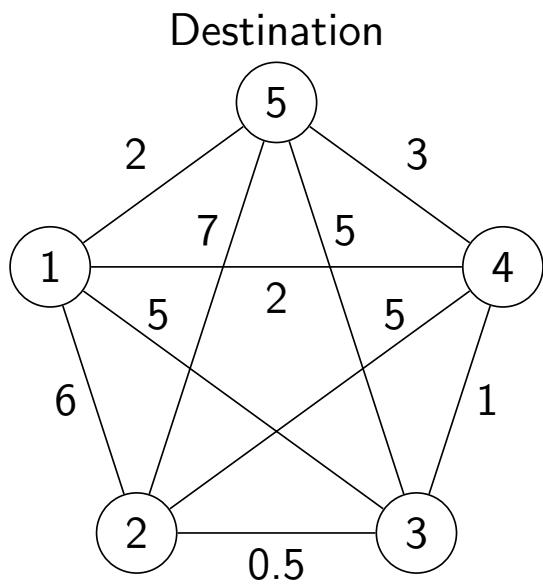
Solving the Shortest Path Problem



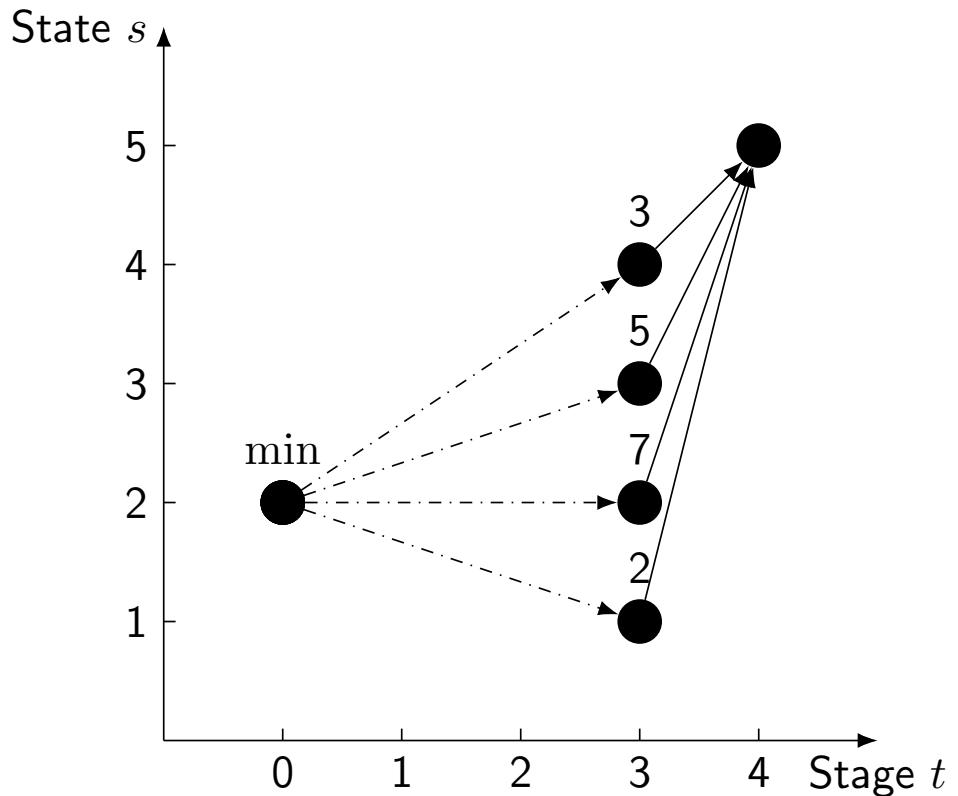
Destination is node 5.



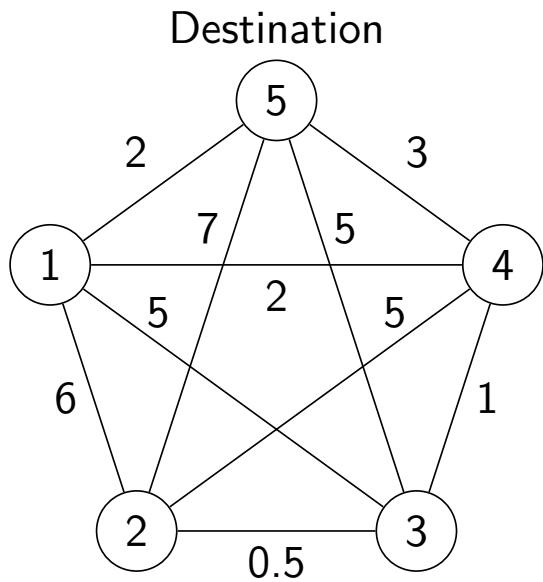
Solving the Shortest Path Problem



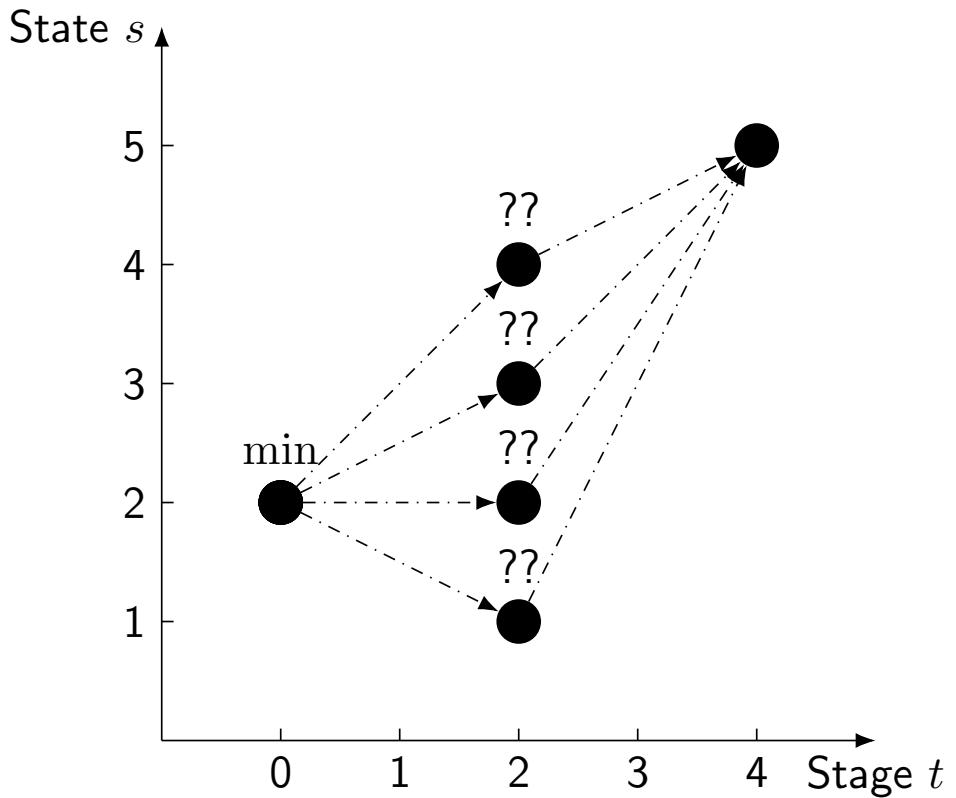
Destination is node 5.



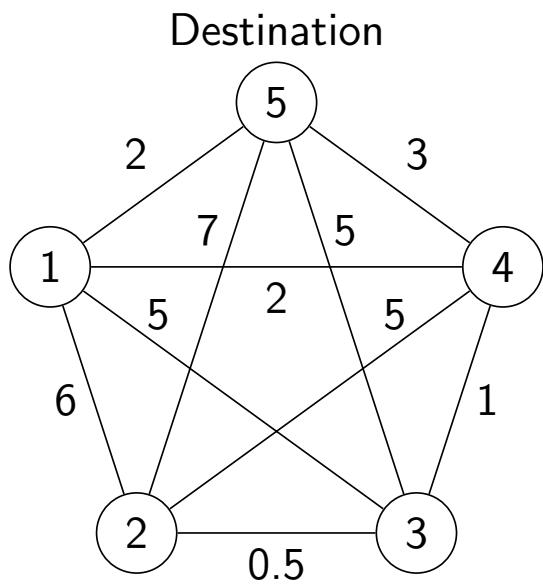
Solving the Shortest Path Problem



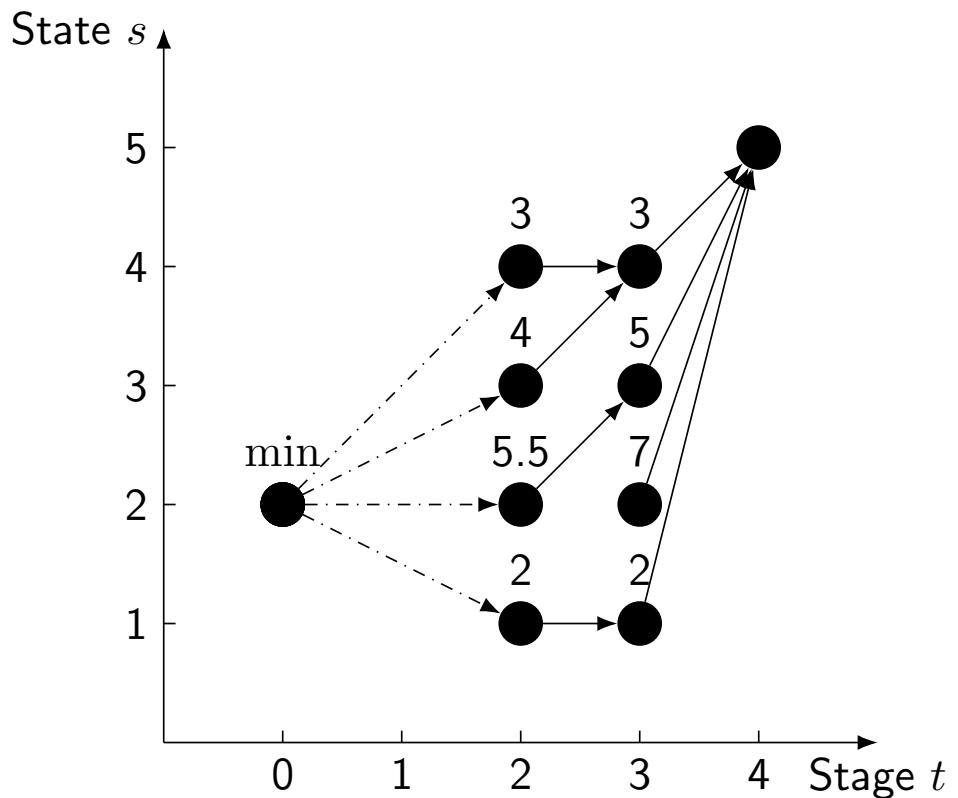
Destination is node 5.



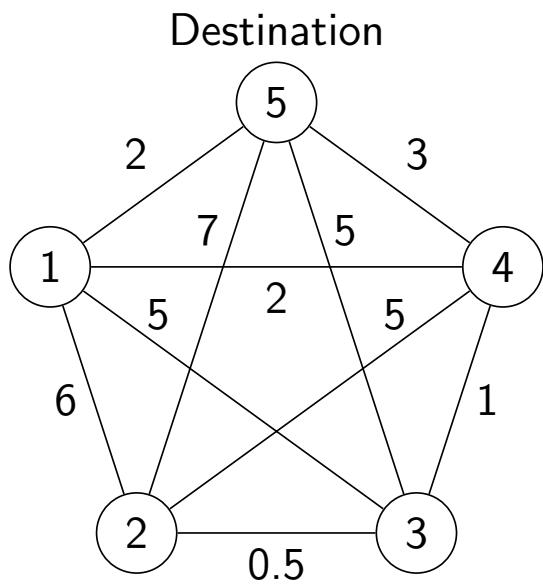
Solving the Shortest Path Problem



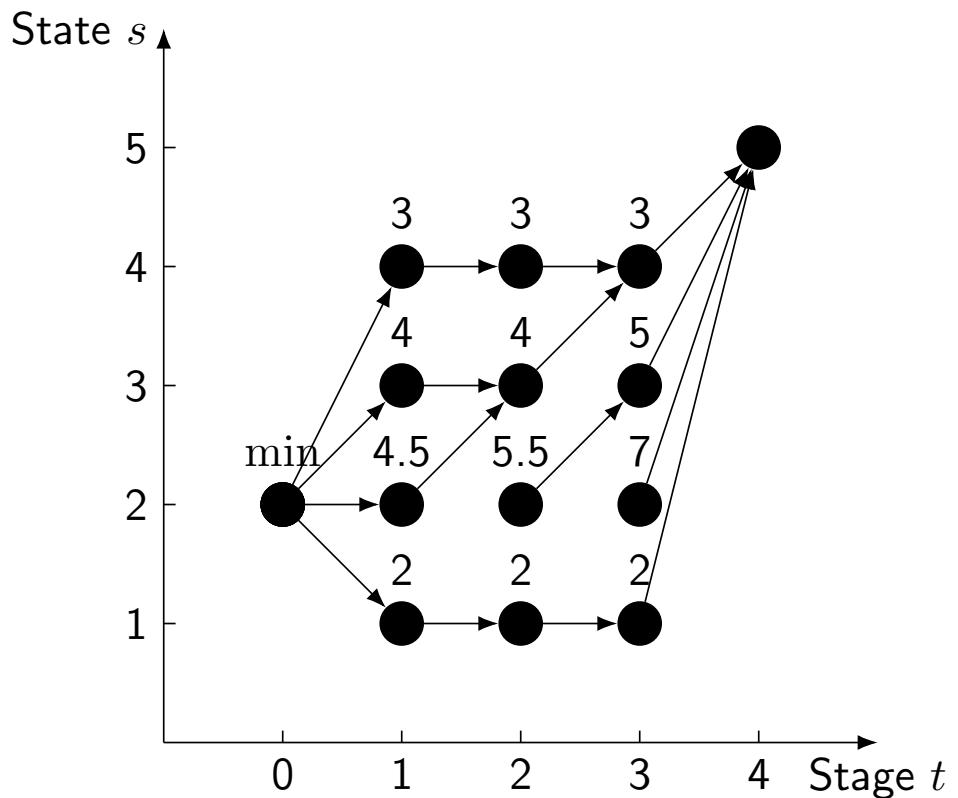
Destination is node 5.



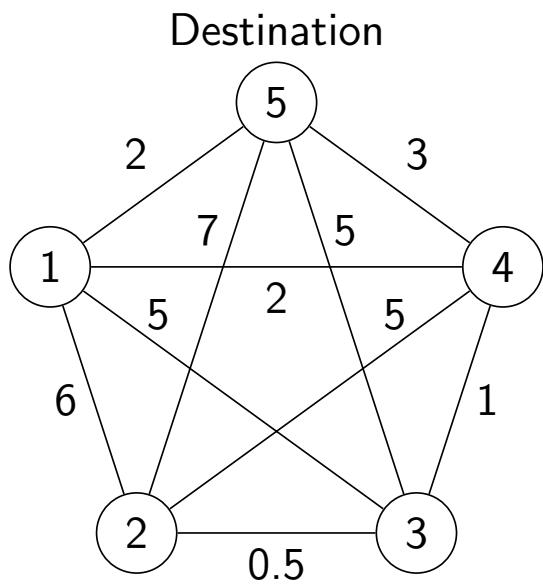
Solving the Shortest Path Problem



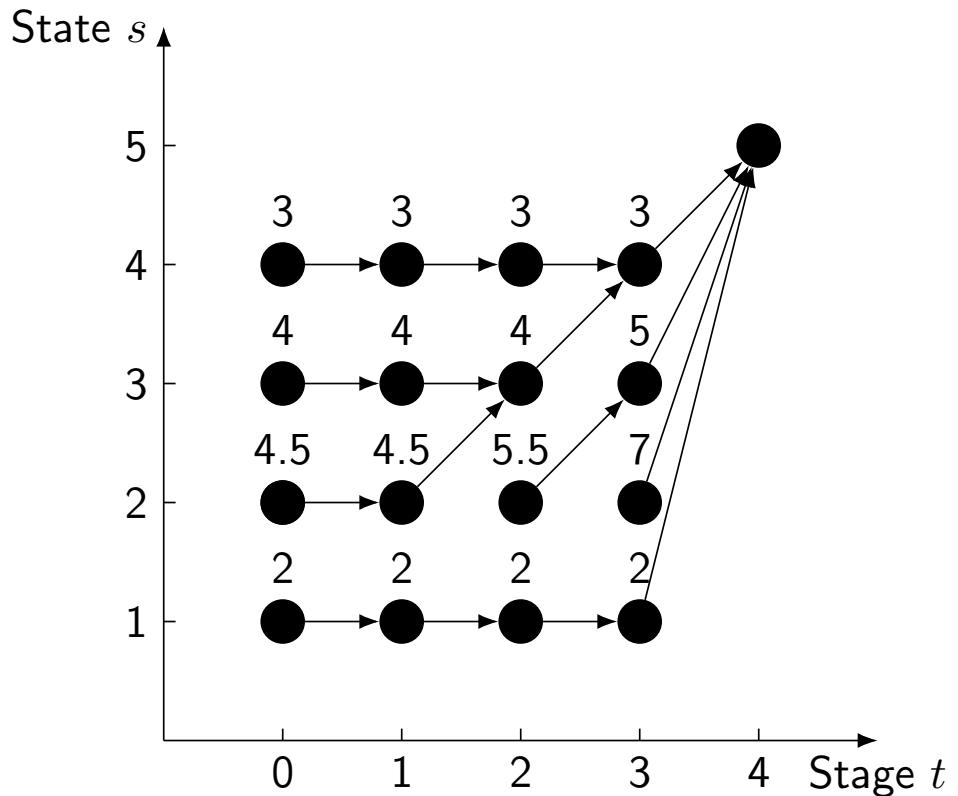
Destination is node 5.



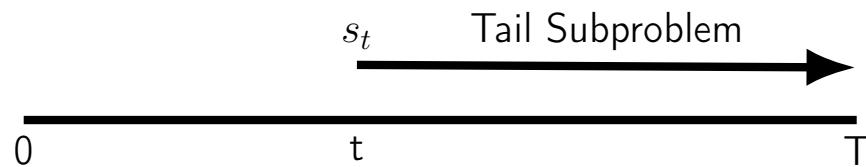
Solving the Shortest Path Problem



Destination is node 5.



Principal of Optimality (Bellman 1957)

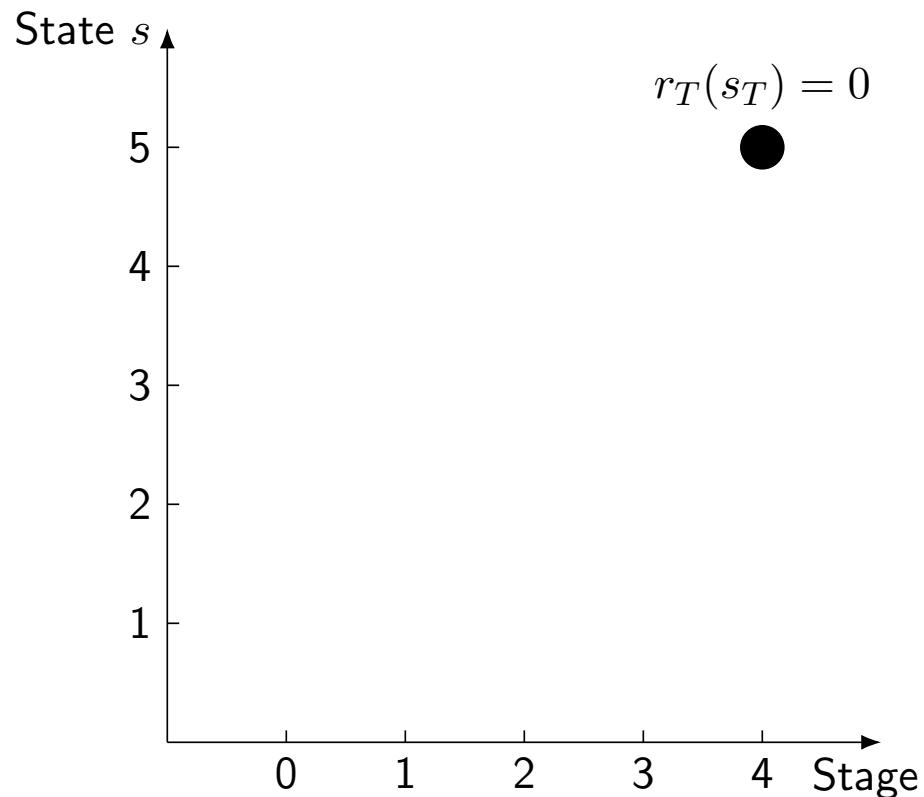


Using Dynamic Programming

$$V_T(s_T) = r_T(s_T)$$

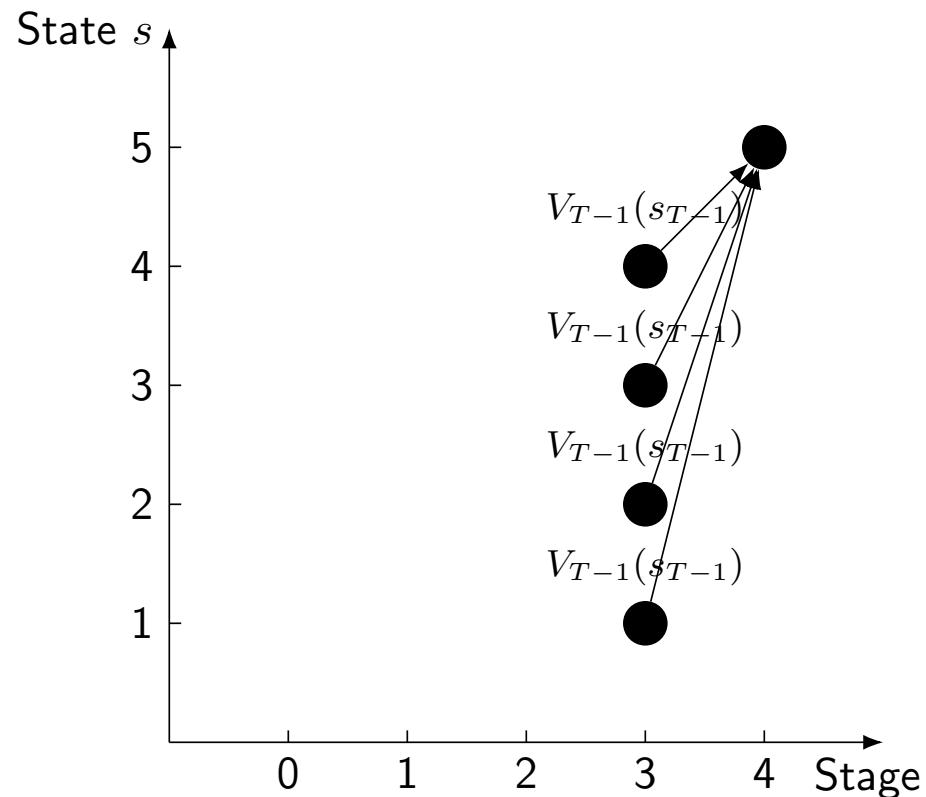
Using Dynamic Programming

```
 $V_T(s_T) = r_T(s_T)$   
for  $t = T - 1, \dots, 0$  do
```



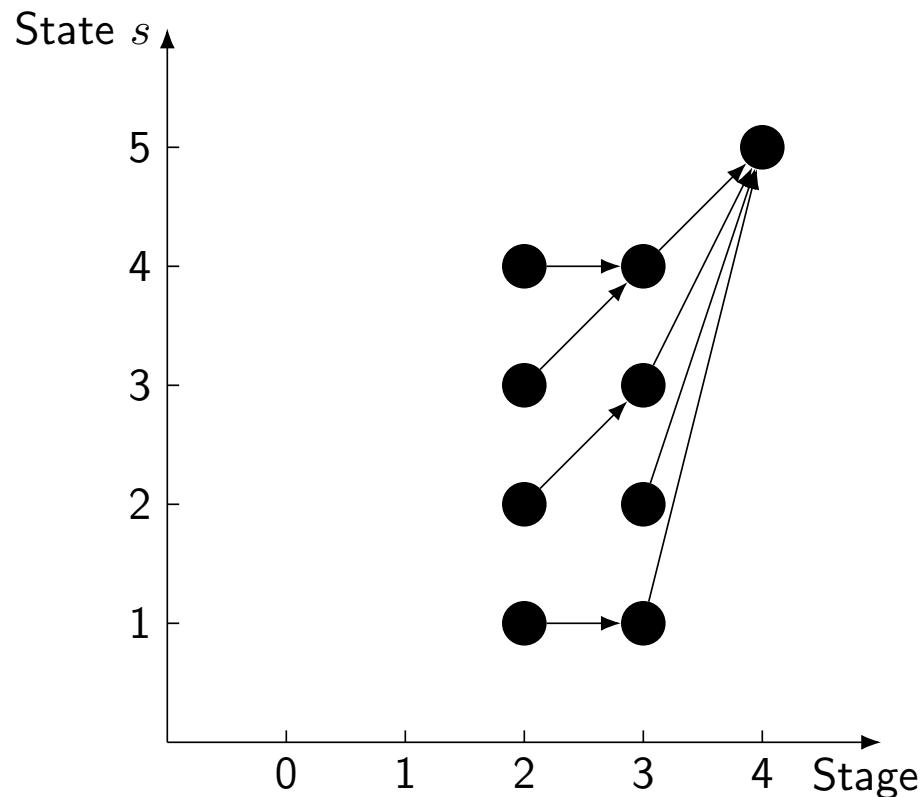
Using Dynamic Programming

```
 $V_T(s_T) = r_T(s_T)$ 
for  $t = T - 1, \dots, 0$  do
     $V_t(s_t) = \max_{a_t \in \mathcal{A}} r_t(s_t, a_t) + V_{t+1}(s_{t+1})$ 
end for
```



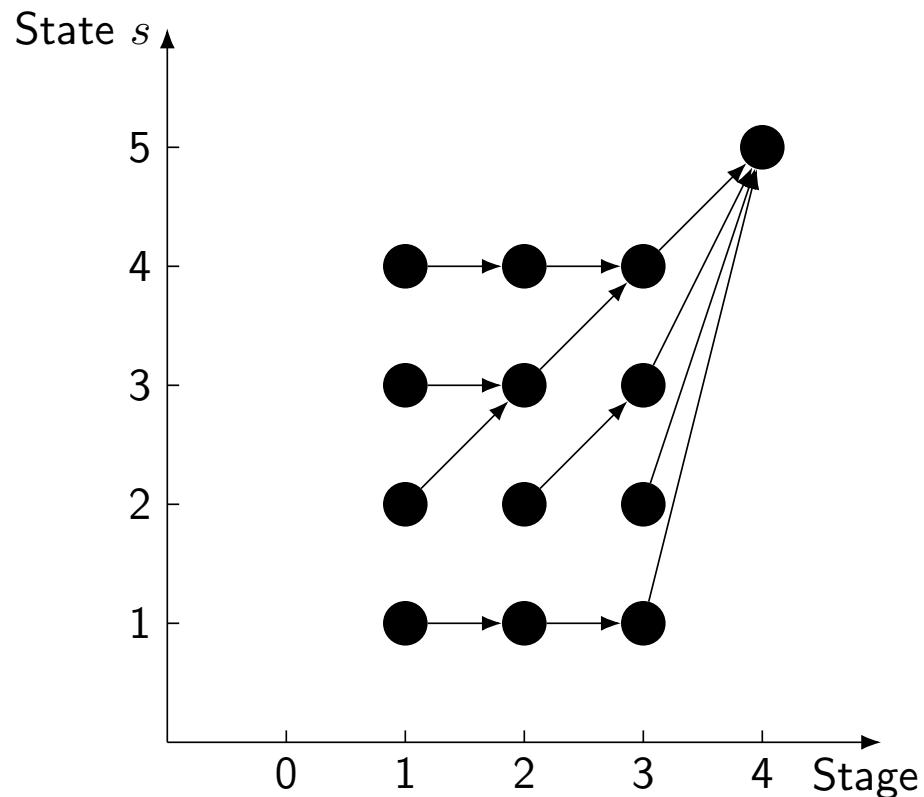
Using Dynamic Programming

```
 $V_T(s_T) = r_T(s_T)$ 
for  $t = T - 1, \dots, 0$  do
     $V_t(s_t) = \max_{a_t \in \mathcal{A}} r_t(s_t, a_t) + V_{t+1}(s_{t+1})$ 
end for
```



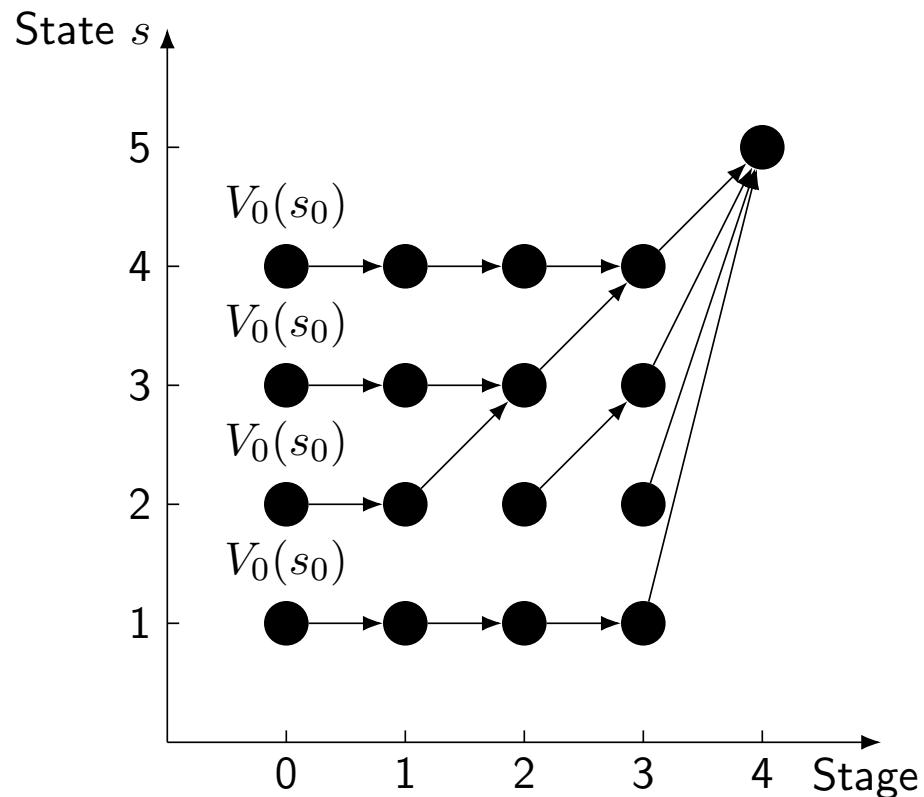
Using Dynamic Programming

```
 $V_T(s_T) = r_T(s_T)$ 
for  $t = T - 1, \dots, 0$  do
     $V_t(s_t) = \max_{a_t \in \mathcal{A}} r_t(s_t, a_t) + V_{t+1}(s_{t+1})$ 
end for
```



Using Dynamic Programming

```
 $V_T(s_T) = r_T(s_T)$ 
for  $t = T - 1, \dots, 0$  do
     $V_t(s_t) = \max_{a_t \in \mathcal{A}} r_t(s_t, a_t) + V_{t+1}(s_{t+1})$ 
end for
```

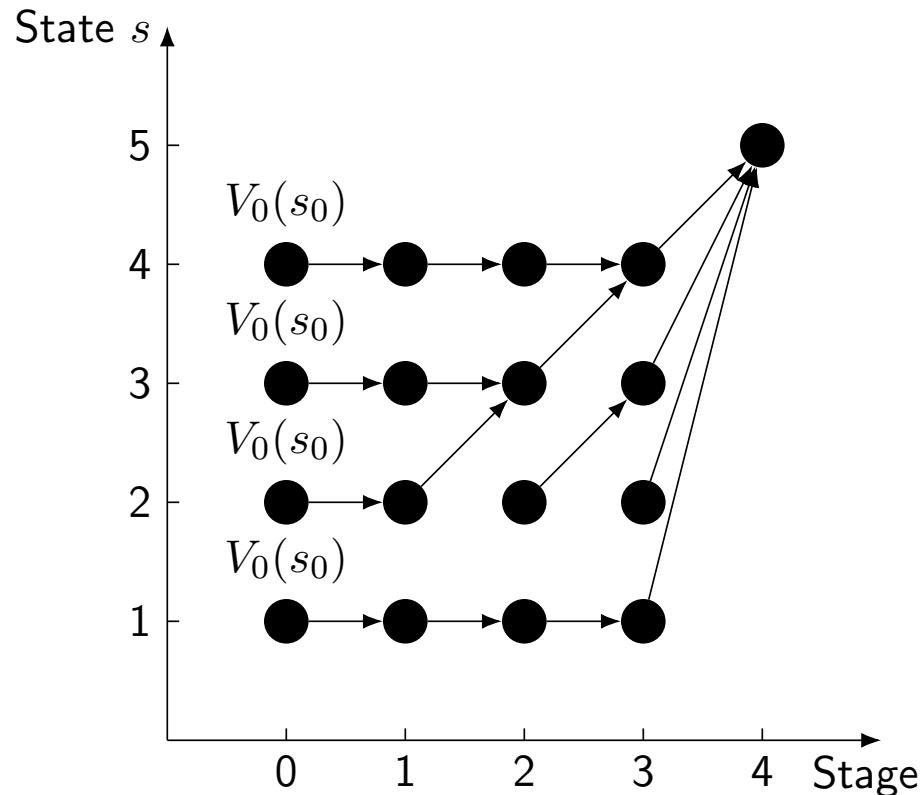


Using Dynamic Programming

```
 $V_T(s_T) = r_T(s_T)$ 
for  $t = T - 1, \dots, 0$  do
     $V_t(s_t) = \max_{a_t \in \mathcal{A}} r_t(s_t, a_t) + V_{t+1}(s_{t+1})$ 
end for
```

Theorem (Dynamic programming)

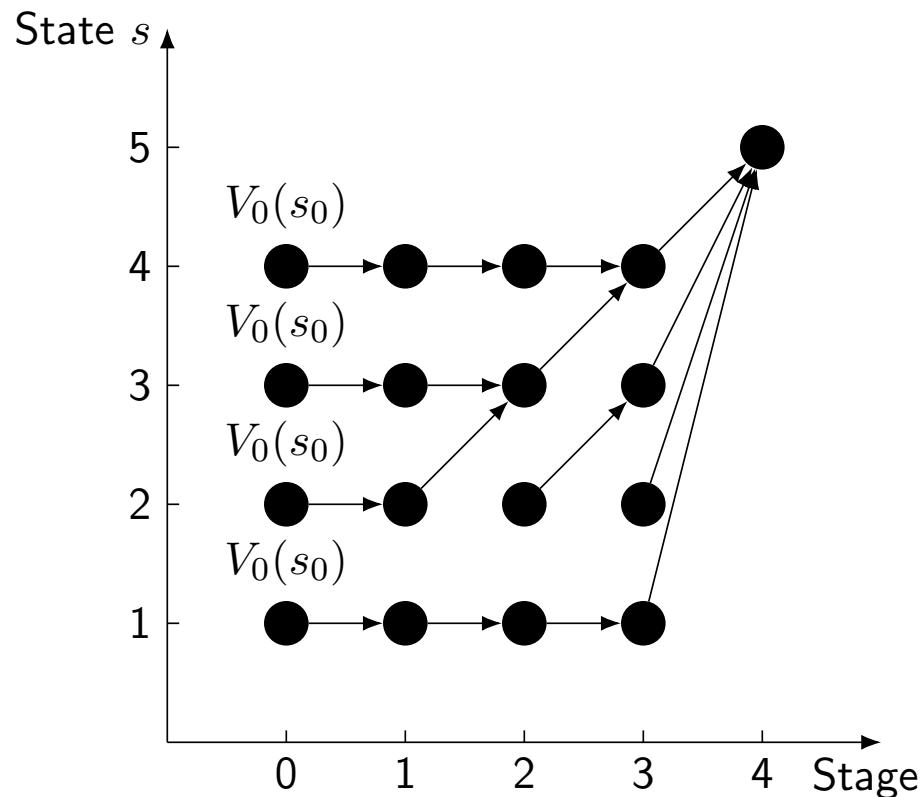
For every initial state s_0 , **the optimal value $V^*(s_0)$ is equal to $V_0(s_0)$** , given above.



Using Dynamic Programming

```
 $V_T(s_T) = r_T(s_T)$ 
for  $t = T - 1, \dots, 0$  do
     $V_t(s_t) = \max_{a_t \in \mathcal{A}} r_t(s_t, a_t) + V_{t+1}(s_{t+1})$ 
end for
```

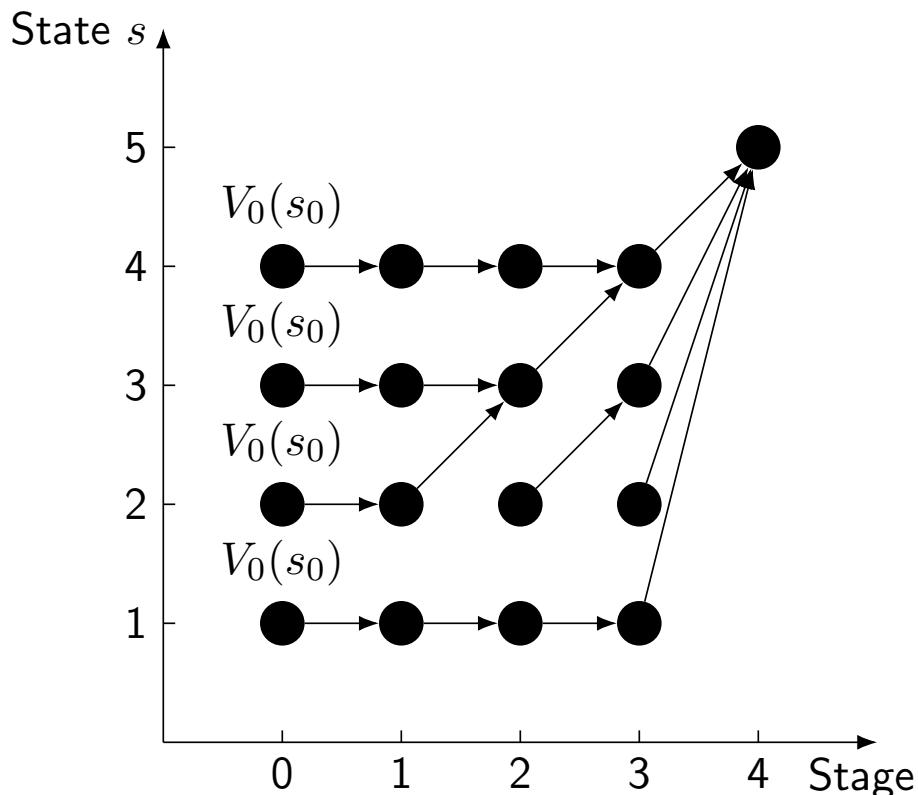
- **Proof:** by induction
- Equivalent to **Bellman-Ford algorithm**
- **Strength:** Generality
- **Weakness:** Computationally expensive



Using Dynamic Programming

```
 $V_T(s_T) = r_T(s_T)$ 
for  $t = T - 1, \dots, 0$  do
     $V_t(s_t) = \max_{a_t \in \mathcal{A}} r_t(s_t, a_t) + V_{t+1}(s_{t+1})$ 
end for
```

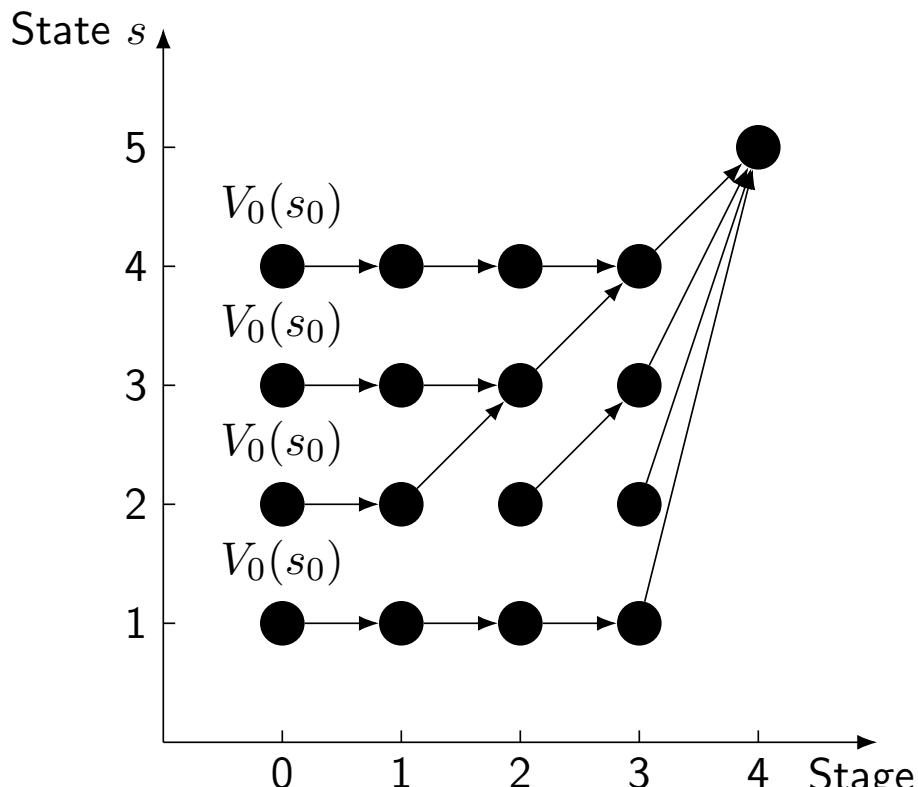
- **Proof:** by induction
- Equivalent to **Bellman-Ford algorithm**
- **Strength:** Generality
- **Weakness:** Computationally expensive
 $\mathcal{O}(|S||A|T)$
- Much better than naive approach $\mathcal{O}(T!)$



Using Dynamic Programming

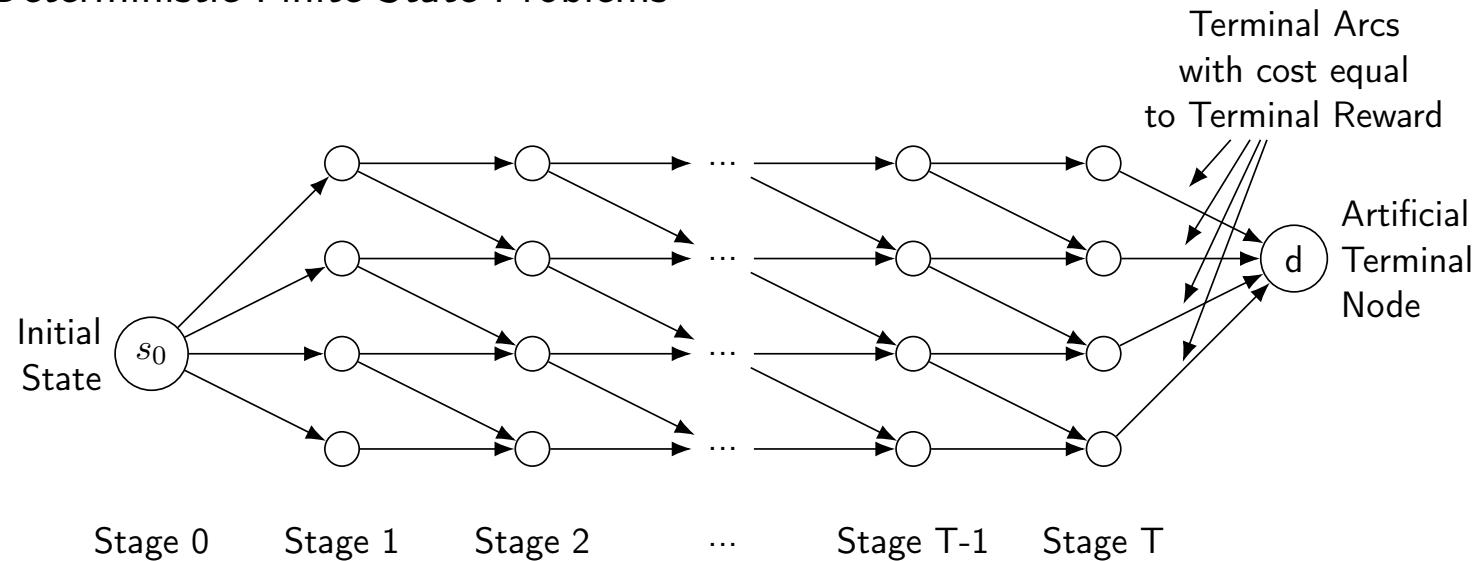
```
 $V_T(s_T) = r_T(s_T)$ 
for  $t = T - 1, \dots, 0$  do
     $V_t(s_t) = \max_{a_t \in \mathcal{A}} r_t(s_t, a_t) + V_{t+1}(s_{t+1})$ 
end for
```

- **Proof:** by induction
- Equivalent to **Bellman-Ford algorithm**
- **Strength:** Generality
- **Weakness:** Computationally expensive
 $\mathcal{O}(|S||A|T)$
- Much better than naive approach $\mathcal{O}(T!)$
- ↗ ALL the tail subproblems are solved
(in addition to the original problem)



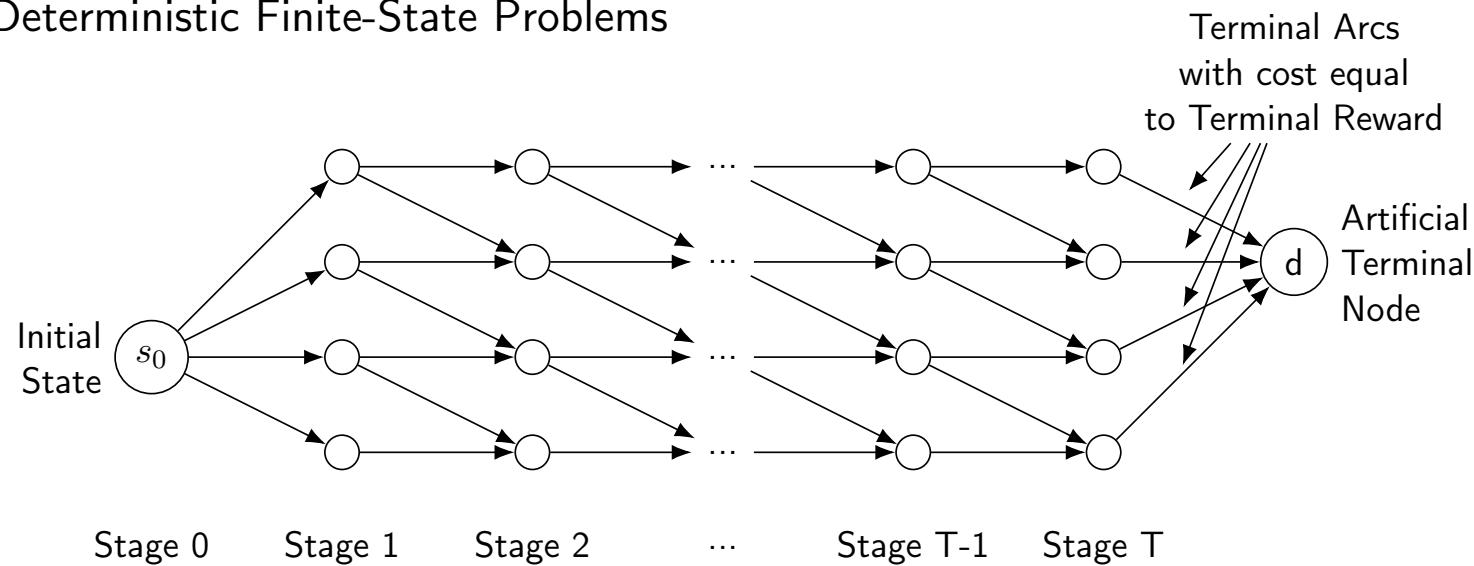
Sequential Decision Making as Shortest Path

For Deterministic Finite-State Problems



Sequential Decision Making as Shortest Path

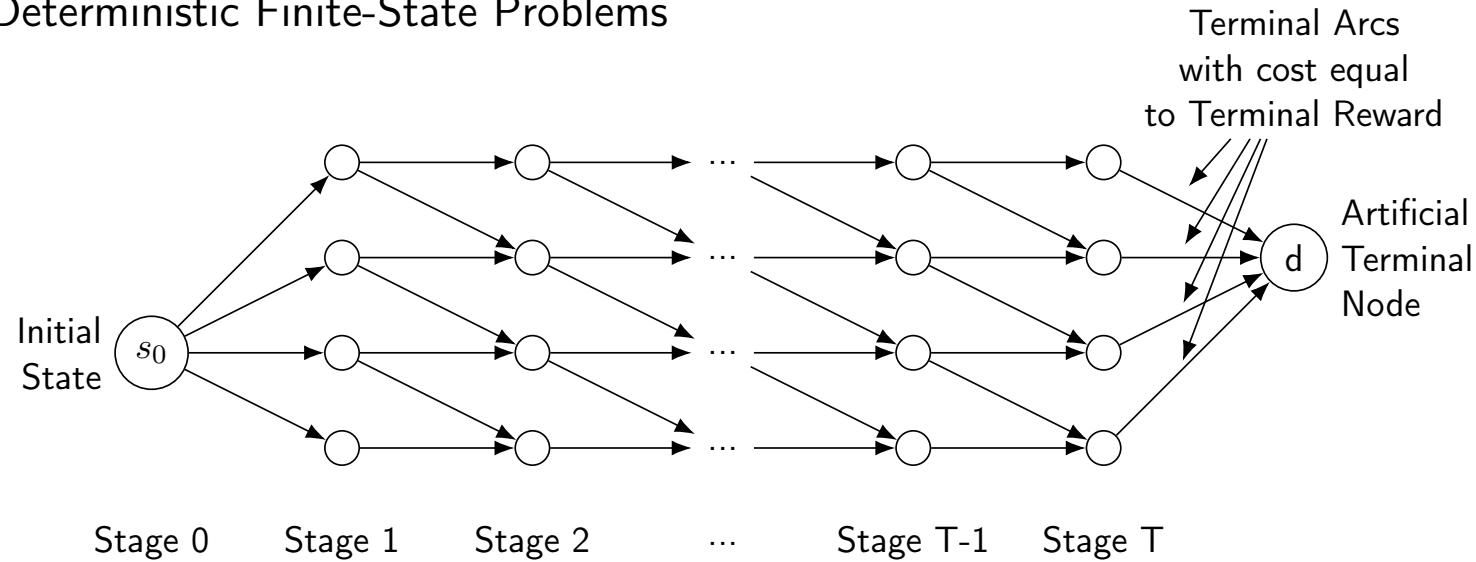
For Deterministic Finite-State Problems



Applications: control systems, industrial manufacturing.

Sequential Decision Making as Shortest Path

For Deterministic Finite-State Problems

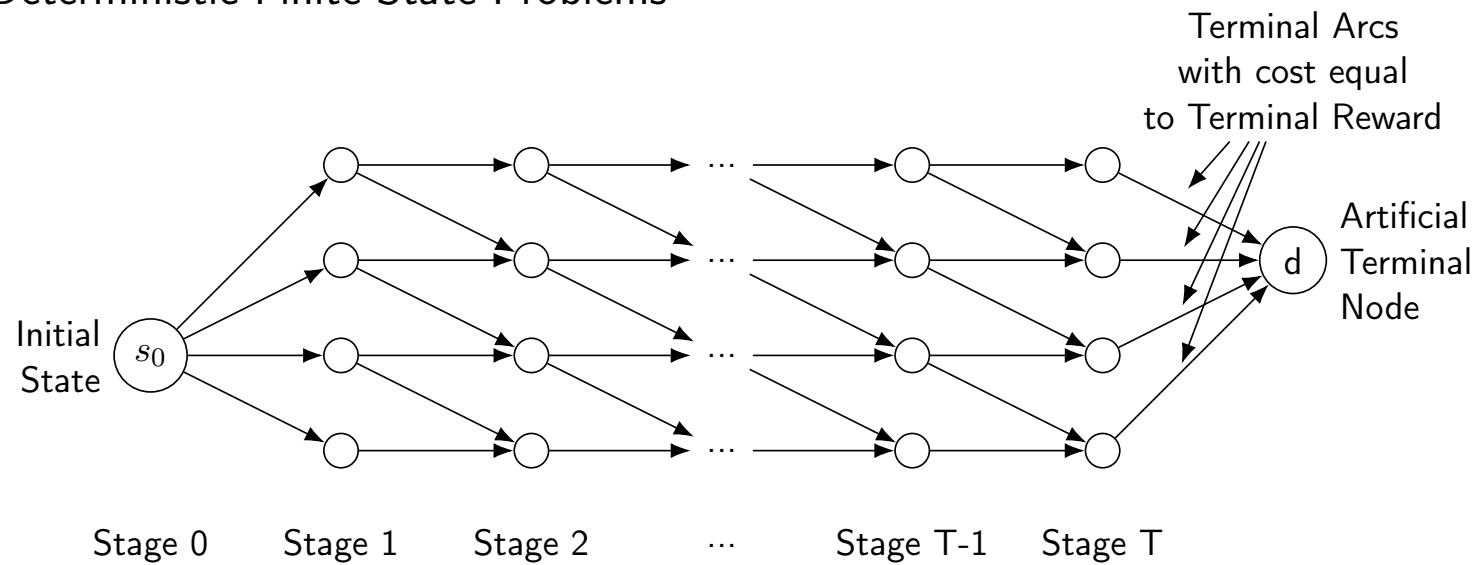


Example: Integer programming (combinatorial optimization)

$$\begin{aligned} & \max \quad c^T x \\ & \text{subject to} \quad Ax = b \\ & \quad x \in \{0, 1\}^T \end{aligned}$$

Sequential Decision Making as Shortest Path

For Deterministic Finite-State Problems



Discuss: If shortest path isn't hard, why are DP problems still challenging?

Sequential Decision Making can get hairy

Example: traveling salesman problem (TSP)

- N cities.
- **Goal:** Find the shortest tour (visit every city exactly once and return home).

Sequential Decision Making can get hairy

Example: traveling salesman problem (TSP)

- N cities.
- **Goal:** Find the shortest tour (visit every city exactly once and return home).
- In this case, can't get around $\mathcal{O}(N!)$. (why?)

Sequential Decision Making can get hairy

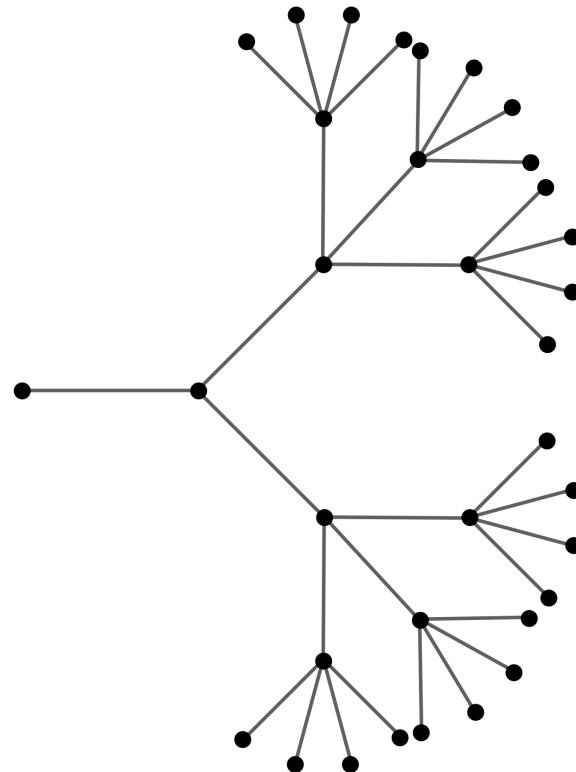
Example: traveling salesman problem (TSP)

- N cities.
- **Goal:** Find the shortest tour (visit every city exactly once and return home).
- In this case, can't get around $\mathcal{O}(N!)$. (why?)
- $|S| = \mathcal{O}(N!)$, $|A| = N$, $T = N$, so $\mathcal{O}(|S||A|T) = \mathcal{O}(N!)$.

Sequential Decision Making can get hairy

Example: traveling salesman problem (TSP)

- N cities.
- **Goal:** Find the shortest tour (visit every city exactly once and return home).
- In this case, can't get around $\mathcal{O}(N!)$. (why?)
- $|S| = \mathcal{O}(N!)$, $|A| = N$, $T = N$, so $\mathcal{O}(|S||A|T) = \mathcal{O}(N!)$.
- This is called the **curse of dimensionality**.



Some Takeaways

- **Sequential decision problems** are those where selected actions affect future states (unlike in bandits).

Some Takeaways

- **Sequential decision problems** are those where selected actions affect future states (unlike in bandits).
- Sequential decision problems are found everywhere. Deterministic examples include routing, combinatorial optimization, linear quadratic control, inventory management.

Some Takeaways

- **Sequential decision problems** are those where selected actions affect future states (unlike in bandits).
- Sequential decision problems are found everywhere. Deterministic examples include routing, combinatorial optimization, linear quadratic control, inventory management.
- The **principle of optimality** relates solving a sequential decision problem to smaller subproblems (called **tail subproblems**).

Some Takeaways

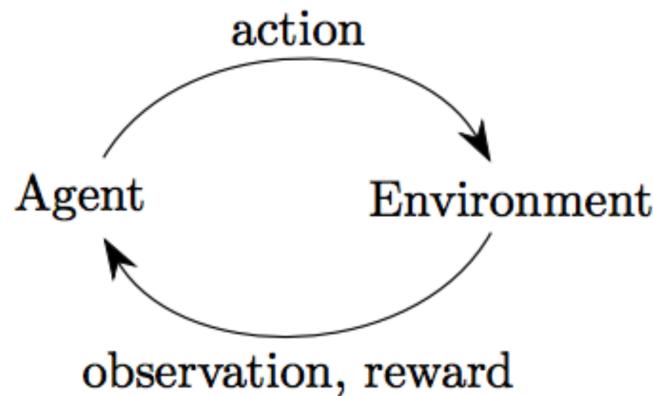
- **Sequential decision problems** are those where selected actions affect future states (unlike in bandits).
- Sequential decision problems are found everywhere. Deterministic examples include routing, combinatorial optimization, linear quadratic control, inventory management.
- The **principle of optimality** relates solving a sequential decision problem to smaller subproblems (called **tail subproblems**).
- **Dynamic programming** solves sequential decision problems by leveraging the principle of optimality.

Some Takeaways

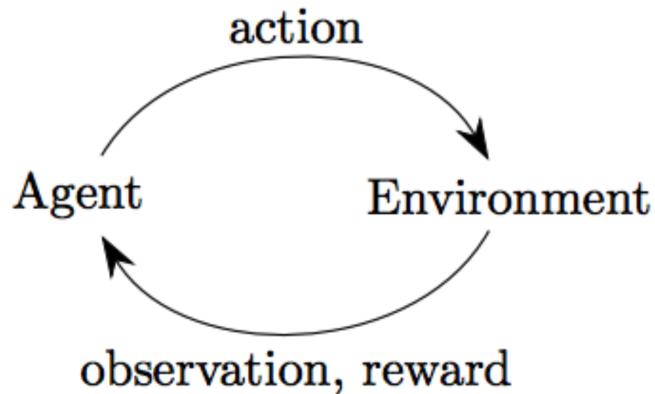
- **Sequential decision problems** are those where selected actions affect future states (unlike in bandits).
- Sequential decision problems are found everywhere. Deterministic examples include routing, combinatorial optimization, linear quadratic control, inventory management.
- The **principle of optimality** relates solving a sequential decision problem to smaller subproblems (called **tail subproblems**).
- **Dynamic programming** solves sequential decision problems by leveraging the principle of optimality.
- The **curse of dimensionality** refers to the exponential growth in state (and action) spaces. This renders “efficient” dynamic programming algorithms insufficient for many problems of interest.

Why Stochastic Problems?

- Stochastic environment:



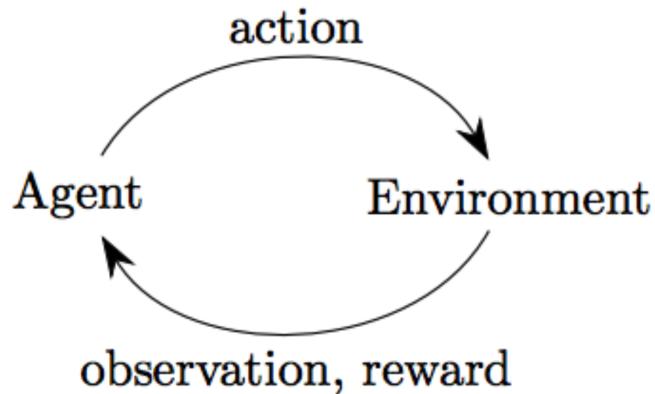
Why Stochastic Problems?



- Stochastic environment:

- Uncertainty in **rewards** (e.g. multi-armed bandits, contextual bandits)

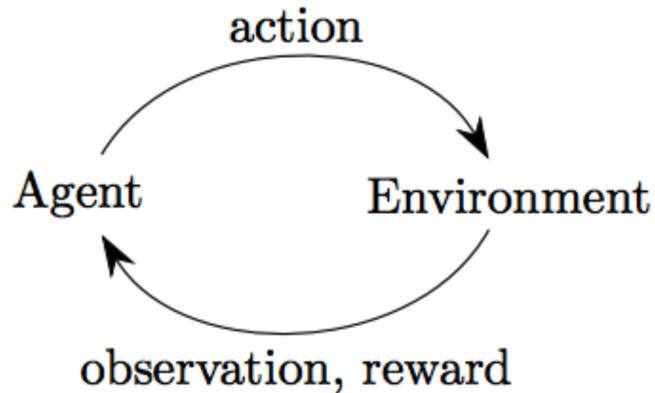
Why Stochastic Problems?



■ Stochastic environment:

- Uncertainty in **rewards** (e.g. multi-armed bandits, contextual bandits)
- Uncertainty in **dynamics**, i.e.
 $(s_t, a_t) \rightarrow s_{t+1}$

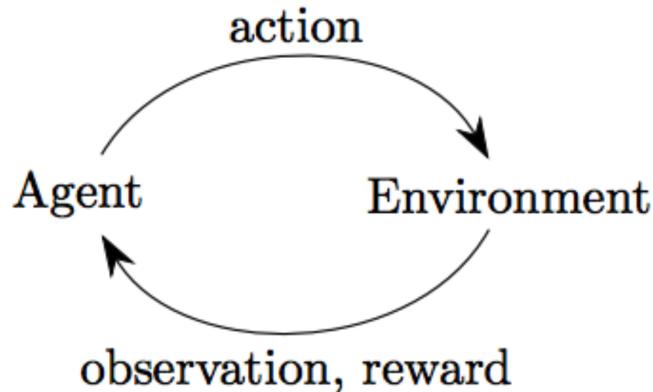
Why Stochastic Problems?



■ Stochastic environment:

- Uncertainty in **rewards** (e.g. multi-armed bandits, contextual bandits)
- Uncertainty in **dynamics**, i.e. $(s_t, a_t) \rightarrow s_{t+1}$
- Uncertainty in **horizon** (called **stochastic shortest path**)

Why Stochastic Problems?

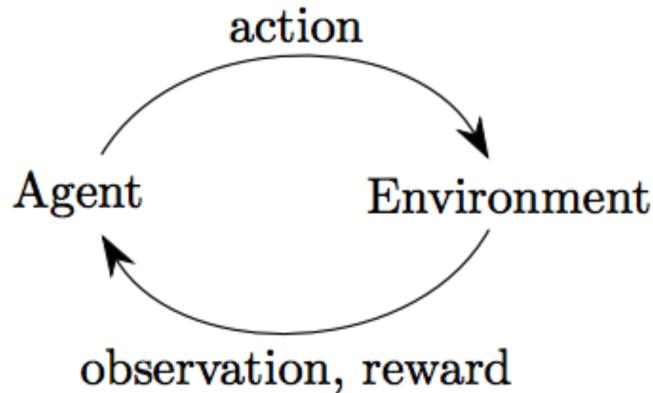


■ Stochastic environment:

- Uncertainty in **rewards** (e.g. multi-armed bandits, contextual bandits)
- Uncertainty in **dynamics**, i.e. $(s_t, a_t) \rightarrow s_{t+1}$
- Uncertainty in **horizon** (called **stochastic shortest path**)

■ Stochastic policies (technical reasons)

Why Stochastic Problems?



■ Stochastic environment:

- Uncertainty in **rewards** (e.g. multi-armed bandits, contextual bandits)
- Uncertainty in **dynamics**, i.e. $(s_t, a_t) \rightarrow s_{t+1}$
- Uncertainty in **horizon** (called **stochastic shortest path**)

■ Stochastic policies (technical reasons)

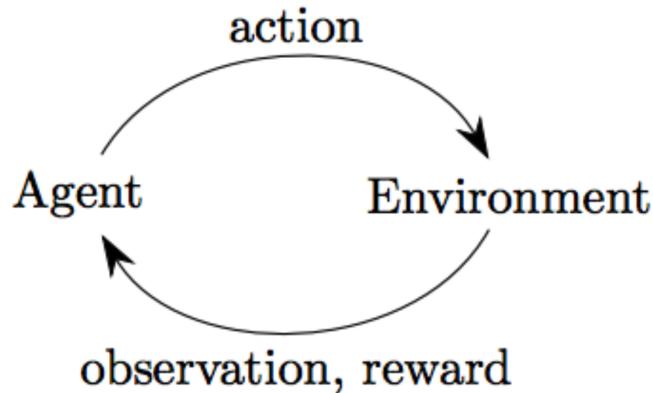
- Trades off **exploration and exploitation**
- Enables **off-policy learning**
- Compatible with **maximum likelihood estimation (MLE)**

Why Stochastic Problems?



- Two-player game of rock–paper–scissors
 - Scissors beats paper
 - Rock beats scissors
 - Paper beats rock
- Consider policies for **iterated** rock–paper–scissors
 - A deterministic policy is **easily exploited**
 - A **uniform random policy** is optimal (i.e., Nash equilibrium)

Why Stochastic Problems?



■ Stochastic environment:

- Uncertainty in **rewards** (e.g. multi-armed bandits, contextual bandits)
- Uncertainty in **dynamics**, i.e. $(s_t, a_t) \rightarrow s_{t+1}$
- Uncertainty in **horizon** (called **stochastic shortest path**)

■ Stochastic policies (technical reasons)

- Trades off **exploration and exploitation**
- Enables **off-policy learning**
- Compatible with **maximum likelihood estimation (MLE)**

👉 Dynamic programming in **deterministic setting** is insufficient.

Example from Supply Chain



Description. At each month t , a warehouse contains s_t items of a specific goods and the demand for that goods is D (stochastic). At the end of each month the manager of the warehouse can order a_t more items from the supplier.

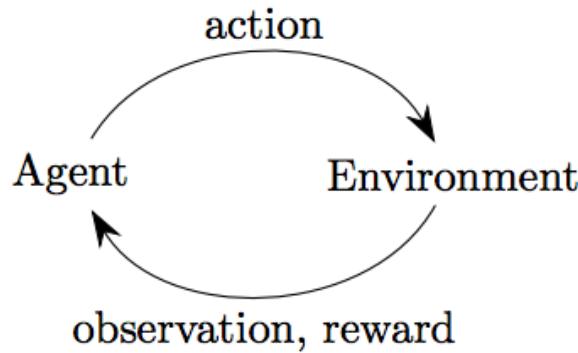
Example from Supply Chain



Description. At each month t , a warehouse contains s_t items of a specific goods and the demand for that goods is D (stochastic). At the end of each month the manager of the warehouse can order a_t more items from the supplier. Furthermore we know that

- The *cost* of maintaining an inventory of s is $h(s)$.
- The *cost* to order a items is $C(a)$.
- The *income* for selling q items is $f(q)$.
- If the demand $d \sim D$ is bigger than the available inventory s , customers that cannot be served leave.
- The *value of the remaining inventory* at the end of the year is $g(s)$.
- *Constraint:* the store has a maximum capacity M .

The problem



Goal

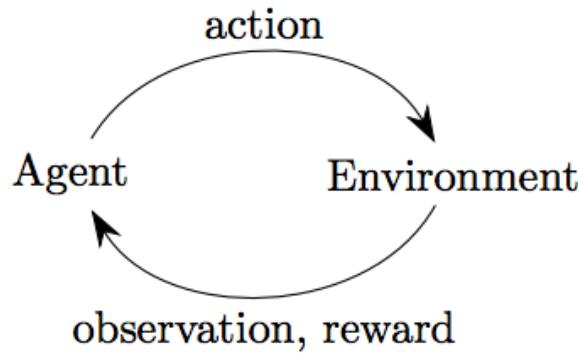
$$a_t = \pi(s_{0:t}; \theta)$$

s₁, a₁, r₁, s₂, a₂, r₂, s₃, a₃, r₃, s₄, a₄, r₄, s₅, a₅, r₅,

(State-action-reward trajectory)

(trajectory or rollout)

The problem



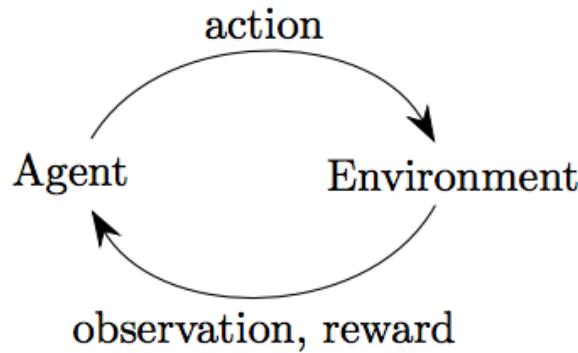
Goal

$$a_t = \pi(s_{0:t}; \theta)$$

s1, a1, r1, s2, a2, r2, s3, a3, r3, s4, a4, r4, s5, a5, r5,

State: {red, green}

The problem



Goal

$$a_t = \pi(s_{0:t}; \theta)$$



a1, r1, s2, a2, r2, s3, a3, r3,



s4, a4, r4,



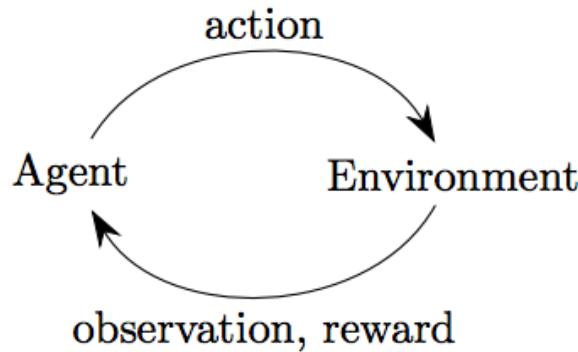
s5, a5, r5,

State: {red, green}

r: 3 times

$$p(a | red)$$

The problem



Goal

$$a_t = \pi(s_{0:t}; \theta)$$

$s_1, a_1, r_1, s_2, a_2, r_2, s_3, a_3, r_3, s_4, a_4, r_4, s_5, a_5, r_5, \dots$

State: {red, green}

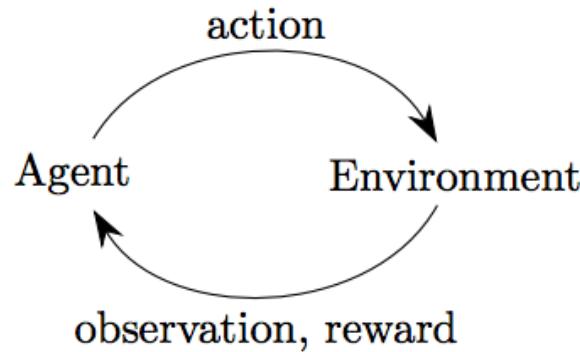
r: 3 times

$$p(a | red)$$

g: 2 times

$$p(a | green)$$

The problem



Goal

$$a_t = \pi(s_{0:t}; \theta)$$

$s_1, a_1, r_1, s_2, a_2, r_2, s_3, a_3, r_3, s_4, a_4, r_4, s_5, a_5, r_5, \dots$

State: {red, green}

r: 3 times

$$p(a | red)$$

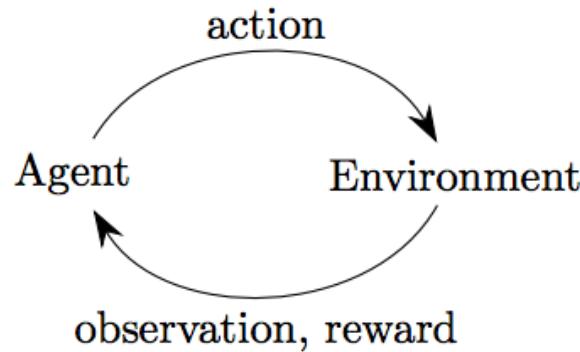
g: 2 times

$$p(a | green)$$

r, g: 1 time

$$p(a | red, green)$$

The problem



Goal

$$a_t = \pi(s_{0:t}; \theta)$$

s1, a1, r1, s2, a2, r2, s3, a3, r3, s4, a4, r4, s5, a5, r5,

State: {red, green}

r: 3 times

g: 2 times

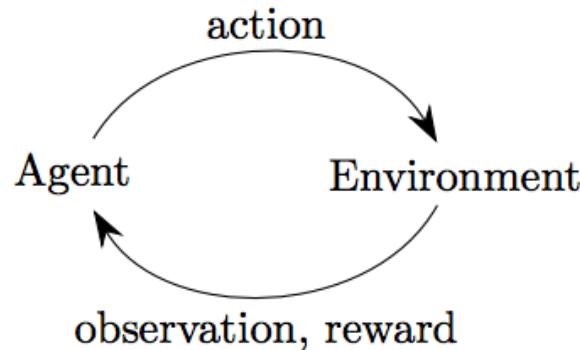
r, g: 1 time

r, g, r: 0 time

#data	N_1	$p(a red)$
	N_2	$p(a green)$
	N_3	$p(a red, green)$
	N_4	$p(a red, green, red)$

$$N_1 \sim N_2 < N_3 << N_4$$

The problem



Goal

$$a_t = \pi(s_{0:t}; \theta)$$

s1, a1, r1, s2, a2, r2, s3, a3, r3, s4, a4, r4, s5, a5, r5,

State: {red, green}

r: 3 times

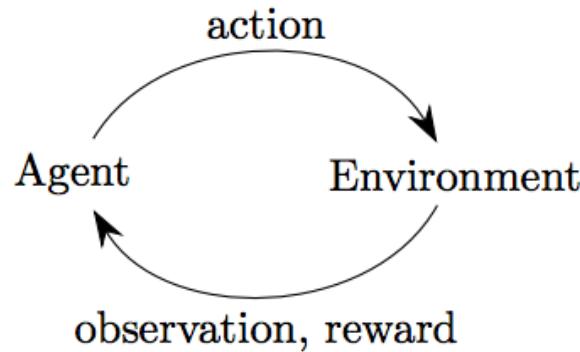
g: 2 times

r, g: 1 time

r, g, r: 0 time

r, r, r, r, r, r, r, g, r: very unlikely!

The problem



Goal

$$a_t = \pi(s_{0:t}; \theta)$$

s1, a1, r1, s2, a2, r2, s3, a3, r3, s4, a4, r4, s5, a5, r5,

State: {red, green}

r: 3 times

g: 2 times

r, g: 1 time

r, g, r: 0 time

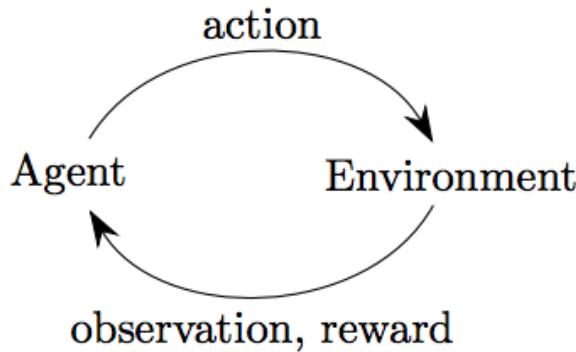
$$p(a | r, \dots, r, g, r)$$

r, r, r, r, r, r, r, g, r: very unlikely!

$$a_t \longleftrightarrow s_{0:t}$$

Data sparsity problem
increases
with t!

The problem



Goal

$$a_t = \pi(s_{0:t}; \theta)$$

Learning issues!

s1, a1, r1, s2, a2, r2, s3, a3, r3, s4, a4, r4, s5, a5, r5,

State: {red, green}

r: 3 times

g: 2 times

r, g: 1 time

r, g, r: 0 time

$$p(a | r, \dots, r, g, r)$$

r, r, r, r, r, r, r, g, r: very unlikely!

$$a_t \longleftrightarrow s_{0:t}$$

Data sparsity problem
increases
with t!

Problem Formulation

$$\max \sum_{t=1}^T r_t$$

environment model

policy

s.t. $s_{t+1} = f(s_{0:t}, a_{0:t}) \quad a_t = \pi(s_{0:t}; \theta)$



$$s_{t+1} = f(s_t, a_t) \quad a_t = \pi(s_t; \theta)$$

markov assumption

Problem Formulation

$$\max \sum_{t=1}^T r_t$$

environment model

policy

s.t. $s_{t+1} = f(s_t, a_t)$ $a_t = \pi(s_t; \theta)$

$$r_t = g(s_t, a_t)$$

markov reward function

markov decision process (MDP)

MDP Formulation

$$\max \sum_{t=1}^T r_t$$

s.t. $s_{t+1} = f(s_t, a_t)$ $a_t = \pi(s_t; \theta)$



s_t :(position, velocity)

is this sufficient?

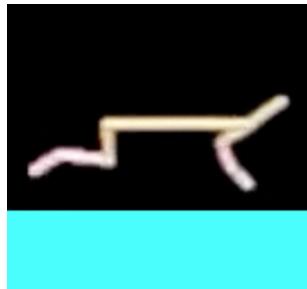
$$r_t = g(s_t, a_t)$$

$$s_{t+1} = f(s_t, s_{t-1}, a_t)$$
 Yes!

MDP Formulation

$$\max \sum_{t=1}^T r_t$$

s.t. $s_{t+1} = f(s_t, a_t)$ $a_t = \pi(s_t; \theta)$



s_t :(position, velocity)

$$r_t = g(s_t, a_t)$$

is this sufficient?

s_t :(position, velocity, acceleration)

Always expand the state-space to make any problem Markov!

comes at a cost ... data sparsity ...

Policy Optimization

$$\tau = (s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_{t-1}, a_{t-1}, r_{t-1}, s_t)$$

Trajectory

$$\max_{\theta} E_{\tau}[R(\tau)]$$

policy $a_t = \pi(s_t; \theta)$

Consider the Value Function

Value Function: Expected Sum of Rewards

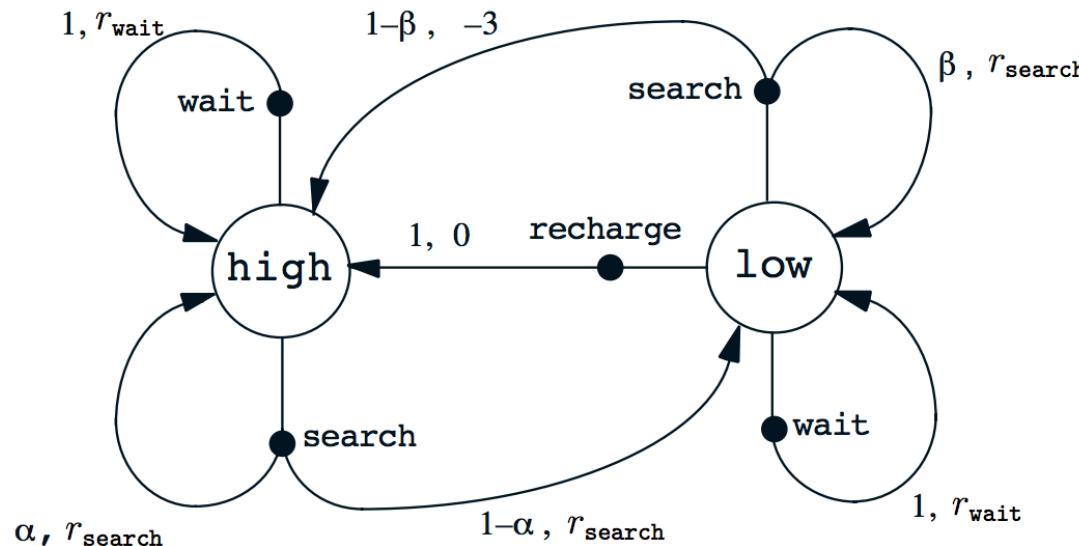
$$V(s_{t_0}) = E_\tau \left[\sum_{t=t_0} \gamma^{t-t_0} r_t \right]$$

$\gamma < 1$: discount factor

$$V(s_{t_0}) = r_{t_0} + \gamma V(s_{t_0+1})$$

Tabular Decision Process Example: Cleaning Robot

Finite State Space



Discrete

s	a	s'	$p(s' s, a)$	$r(s, a, s')$
high	search	high	α	r_{search}
high	search	low	$1 - \alpha$	r_{search}
low	search	high	$1 - \beta$	-3
low	search	low	β	r_{search}
high	wait	high	1	r_{wait}
high	wait	low	0	r_{wait}
low	wait	high	0	r_{wait}
low	wait	low	1	r_{wait}
low	recharge	high	1	0
low	recharge	low	0	0.

Model

Value Function in Tabular Setup

$$V(s_{t_0}) = r_{t_0} + \gamma V(s_{t_0+1})$$

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s' \in S} p(s' | s, \pi(s)) V^\pi(s')$$

Value Function
is for a policy

Value Function in Tabular Setup

$$V(s_{t_0}) = r_{t_0} + \gamma V(s_{t_0+1})$$

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s' \in S} p(s' | s, \pi(s)) V^\pi(s')$$

s' is next state

Value Function in Tabular Setup

$$V(s_{t_0}) = r_{t_0} + \gamma V(s_{t_0+1})$$

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s' \in S} p(s' | s, \pi(s)) V^\pi(s')$$

Immediate reward

discounted future reward

Value Function in Tabular Setup

$$V(s_{t_0}) = r_{t_0} + \gamma V(s_{t_0+1})$$

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s' \in S} p(s' | s, \pi(s)) V^\pi(s')$$

In simplified notation

$$V^\pi(s) = R(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V^\pi(s')$$

Optimal Policy

$$V^\pi(s) = R(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V^\pi(s')$$

Optimal policy $\pi^* : S \rightarrow A$

$$V^{\pi^*}(s) \geq V^\pi(s) \quad \forall s \in S, \pi : S \rightarrow A$$

Does it exist?

How do we find it?

Optimal Policy

Theorem: Optimal policy exists and its value function satisfies

$$V^\pi(s) = \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V^\pi(s') \right]$$

Bellman Equation

Proof

By Definition,

$$V^\pi(s) = R(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V^\pi(s')$$

Consider, “improvement” of value function

$$V^{\pi'}(s) = \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V^\pi(s') \right]$$

Then

$$V^{\pi'}(s) \geq V^\pi(s), \forall s \in S$$

$$V^{\pi'}(s) > V^\pi(s),$$

if $V^\pi(s)$ doesn't satisfy Bellman Eq for s

iterate until convergence

if convergence to a unique value then, $V^{\pi^*}(s) = V^{\pi'}(s) \forall s$

Convergence to Unique Value

At iteration $k \geq 1$

$$V^{\pi^{k+1}}(s) = \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V^{\pi^k}(s') \right]$$

Difference from optimum

$$\|V^{\pi^k} - V^{\pi^*}\|_\infty = \max_{s \in S} |V^{\pi^k}(s) - V^{\pi^*}(s)|$$

Convergence to Unique Value

$$V^{\pi^{k+1}}(s) - V^{\pi^*}(s)$$

$$= \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V^{\pi^k}(s') \right] - \max_{a' \in A} \left[R(s, a') + \gamma \sum_{s' \in S} p(s'|s, a') V^{\pi^*}(s') \right]$$

$$= \max_{a \in A} \left[\left[R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V^{\pi^k}(s') \right] - \max_{a' \in A} \left[R(s, a') + \gamma \sum_{s' \in S} p(s'|s, a') V^{\pi^*}(s') \right] \right]$$

$$\leq \max_{a \in A} \left[\left[R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V^{\pi^k}(s') \right] - \left[R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V^{\pi^*}(s') \right] \right]$$

Convergence to Unique Value

$$V^{\pi^{k+1}}(s) - V^{\pi^*}(s)$$

$$= \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V^{\pi^k}(s') \right] - \max_{a' \in A} \left[R(s, a') + \gamma \sum_{s' \in S} p(s'|s, a') V^{\pi^*}(s') \right]$$

$$= \max_{a \in A} \left[\left[R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V^{\pi^k}(s') \right] - \max_{a' \in A} \left[R(s, a') + \gamma \sum_{s' \in S} p(s'|s, a') V^{\pi^*}(s') \right] \right]$$

$$\leq \max_{a \in A} \left[\left[\cancel{R(s, a)} + \gamma \sum_{s' \in S} p(s'|s, a) V^{\pi^k}(s') \right] - \left[\cancel{R(s, a)} + \gamma \sum_{s' \in S} p(s'|s, a) V^{\pi^*}(s') \right] \right]$$

Convergence to Unique Value

$$V^{\pi^{k+1}}(s) - V^{\pi^*}(s)$$

$$\leq \gamma \max_{a \in A} \left[\sum_{s' \in S} p(s'|s, a) |V^{\pi^k}(s') - V^{\pi^*}(s')| \right]$$

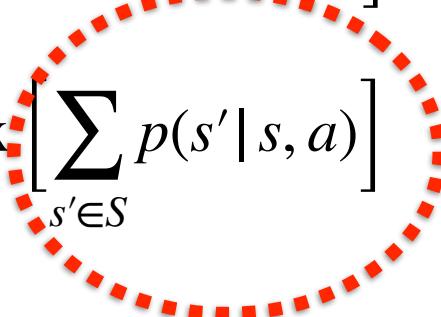
$$\leq \max_{a \in A} \left[\left[R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V^{\pi^k}(s') \right] - \left[R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a) V^{\pi^*}(s') \right] \right]$$

Convergence to Unique Value

$$V^{\pi^{k+1}}(s) - V^{\pi^*}(s)$$

$$\leq \gamma \max_{a \in A} \left[\sum_{s' \in S} p(s' | s, a) | V^{\pi^k}(s') - V^{\pi^*}(s') | \right]$$

$$\leq \gamma \max_{a \in A} \left[\sum_{s' \in S} p(s' | s, a) \| V^{\pi^k}(s') - V^{\pi^*}(s') \|_{\infty} \right]$$

$$\begin{aligned} &= \gamma \| V^{\pi^k}(s') - V^{\pi^*}(s') \|_{\infty} \max_{a \in A} \left[\sum_{s' \in S} p(s' | s, a) \right] \\ &= \gamma \| V^{\pi^k}(s') - V^{\pi^*}(s') \|_{\infty} \end{aligned}$$


Similarly,

$$V^{\pi^*}(s') - V^{\pi^{k+1}}(s') \leq \gamma \| V^{\pi^k}(s') - V^{\pi^*}(s') \|_{\infty}$$

Convergence to Unique Value

$$V^{\pi^{k+1}}(s) - V^{\pi^*}(s) = \gamma \|V^{\pi^k}(s') - V^{\pi^*}(s')\|_\infty$$

$$\|V^{\pi^*}(s') - V^{\pi^{k+1}}(s')\|_\infty \leq \gamma \|V^{\pi^k}(s') - V^{\pi^*}(s')\|_\infty$$

Similarly,

$$V^{\pi^*}(s') - V^{\pi^{k+1}}(s') \leq \gamma \|V^{\pi^k}(s') - V^{\pi^*}(s')\|_\infty$$

Convergence to Unique Value

$$\|V^{\pi^*}(s') - V^{\pi^{k+1}}(s')\|_{\infty} \leq \gamma \|V^{\pi^k}(s') - V^{\pi^*}(s')\|_{\infty}$$

Since $\gamma < 1$

$$\|V^{\pi^*}(s') - V^{\pi^{k+1}}(s')\|_{\infty} \rightarrow 0$$

Optimal Value Function is Unique

Value Iteration

$$V^{\pi^{k+1}}(s) = \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V^{\pi^k}(s') \right]$$

Value iteration converges to optimum value

$$\pi^*(s) = \arg \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V^{\pi^*}(s') \right]$$

In General

for N steps

Policy Evaluation

$$V^{\pi^{k+1}}(s) = R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a)V^{\pi^k}(s')$$

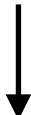
Policy Improvement

$$\pi^*(s) = \arg \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a)V^{\pi^*}(s') \right]$$

Generalized Policy Iteration

Initialize Policy

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$



Evaluate Policy

$$\Delta \leftarrow 0$$

For each $s \in \mathcal{S}$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

Improve Policy

policy-stable \leftarrow true

For each $s \in \mathcal{S}$:

$$\textit{old-action} \leftarrow \pi(s)$$

$$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

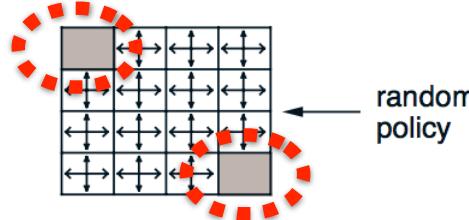
Policy Iteration

$k = 0$

v_k for the
Random Policy

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

Greedy Policy
w.r.t. v_k



High-reward
Locations

Unique Optimal
Value Function

Many Optimal
Policies

Generalized Policy Iteration

for N steps

Policy Evaluation

$$V^{\pi^{k+1}}(s) = R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a)V^{\pi^k}(s')$$

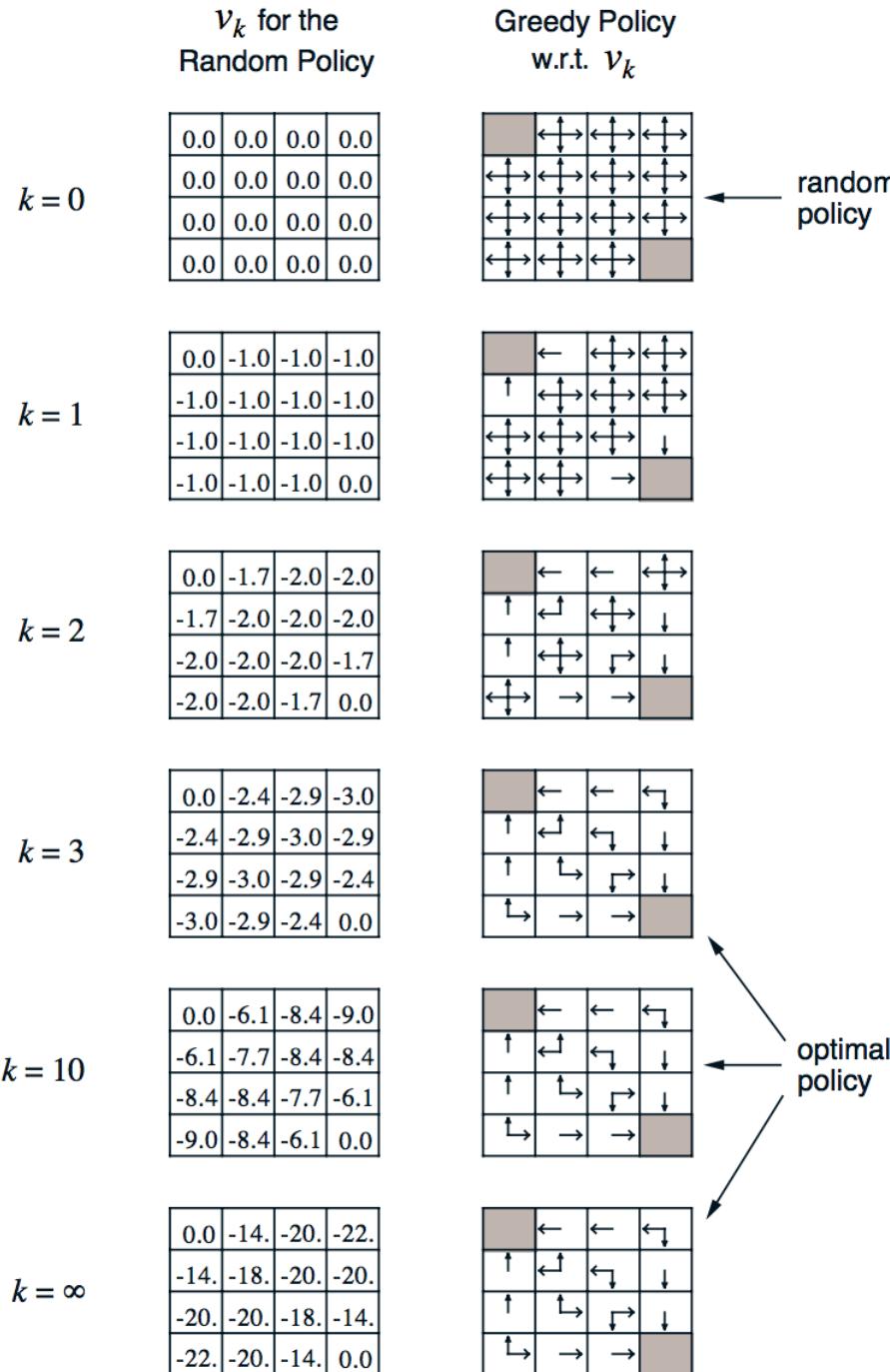
Policy Improvement

$$\pi^*(s) = \arg \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a)V^{\pi^*}(s') \right]$$

Value Iteration (special case N=1)

$$V^{\pi^{k+1}}(s) = \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} p(s'|s, a)V^{\pi^k}(s') \right]$$

Policy iteration
can converge
faster than
value iteration!



Unique Optimal Value Function

Many Optimal Policies

Value Iteration

$$V^{\pi^{k+1}}(s) = \max_{a \in A} \left[R(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) V^{\pi^k}(s') \right]$$

Q-Function

$$Q(s_{t_0}, a_{t_0}) = r(s_{t_0}, a_{t_0}) + \gamma E_\tau \left[\sum_{t=t_0+1} \gamma^{t-(t_0+1)} r_t \right]$$

Q-Value Iteration

$$Q^{\pi^{k+1}}(s, a) = R(s, a) + \gamma \sum_{s' \in S} p(s' | s, a) \max_{a'} Q^{\pi^k}(s', a')$$