# Markov Decision Processes
## Chapter 17

Mausam

# Planning Agent



**Static vs. Dynamic**

**Environment**

**Fully vs. Partially Observable**

**Perfect vs. Noisy**

*What action next?*

**Deterministic vs. Stochastic**

**Instantaneous vs. Durative**

*Percepts*

*Actions*

2

# Classical Planning

**Static**

Environment

**Fully Observable**

**Deterministic**

*What action next?*

**Instantaneous**

**Perfect**

*Percepts*

*Actions*

3

# Stochastic Planning: MDPs

**Static**

*Environment*

**Fully Observable**

**Stochastic**

**What action next?**

**Instantaneous**

**Perfect**

*Percepts*

*Actions*

# MDP vs. Decision Theory

- Decision theory - episodic

- MDP -- sequential

# Markov Decision Process (MDP)

- $\mathcal{S}$: A set of states

# Markov Decision Process (MDP)

- $\mathcal{S}$: A set of states
- $\mathcal{A}$: A set of actions

# Markov Decision Process (MDP)

- $\mathcal{S}$: A set of states
- $\mathcal{A}$: A set of actions
- $\mathcal{T}$(s,a,s'): transition model

# Markov Decision Process (MDP)

- $\mathcal{S}$: A set of states
- $\mathcal{A}$: A set of actions
- $\mathcal{T}$(s,a,s'): transition model
- $\mathcal{C}$(s,a,s'): cost model

# Markov Decision Process (MDP)

- $\mathcal{S}$: A set of states
- $\mathcal{A}$: A set of actions
- $\mathcal{T}(s,a,s')$: transition model
- $\mathcal{C}(s,a,s')$: cost model

# Markov Decision Process (MDP)

- $\mathcal{S}$: A set of states
- $\mathcal{A}$: A set of actions
- $\mathcal{T}$(s,a,s'): transition model
- $\mathcal{C}$(s,a,s'): cost model
- $\mathcal{G}$: set of goals

# Markov Decision Process (MDP)

- $\mathcal{S}$: A set of states
- $\mathcal{A}$: A set of actions
- $\mathcal{T}$(s,a,s'): transition model
- $\mathcal{C}$(s,a,s'): cost model
- $\mathcal{G}$: set of goals
- $s_0$: start state

# Markov Decision Process (MDP)

- $\mathcal{S}$: A set of states
- $\mathcal{A}$: A set of actions
- $\mathcal{T}$(s,a,s'): transition model
- $\mathcal{C}$(s,a,s'): cost model
- $\mathcal{G}$: set of goals
- $s_0$: start state
- $\gamma$: discount factor

# Markov Decision Process (MDP)

- $\mathcal{S}$: A set of states
- $\mathcal{A}$: A set of actions
- $\mathcal{T}(s,a,s')$: transition model
- $\mathcal{C}(s,a,s')$: cost model
- $\mathcal{G}$: set of goals
- $s_0$: start state
- $\gamma$: discount factor
- $\mathcal{R}(s,a,s')$: reward model

# Markov Decision Process (MDP)

- $\mathcal{S}$: A set of states
- $\mathcal{A}$: A set of actions
- $\mathcal{T}(s,a,s')$: transition model


- $s_0$: start state
- $\gamma$: discount factor
- $\mathcal{R}(s,a,s')$: reward model

# Markov Decision Process (MDP)

- $\mathcal{S}$: A set of states
- $\mathcal{A}$: A set of actions
- $\mathcal{T}$(s,a,s'): transition model



- $\gamma$: discount factor
- $\mathcal{R}$(s,a,s'): reward model

# Markov Decision Process (MDP)

- $\mathcal{S}$: A set of states
- $\mathcal{A}$: A set of actions
- $\mathcal{T}$(s,a,s'): transition model
- $\mathcal{C}$(s,a,s'): cost model
- $\mathcal{G}$: set of goals
- $s_0$: start state

# Objective of an MDP

- Find a policy $\pi: \mathcal{S} \rightarrow \mathcal{A}$

- which optimizes
  - minimizes $\left\{\begin{matrix}\text{discounted}\\\text{or}\\\text{undiscount.}\end{matrix}\right\}$ expected cost to reach a goal
  - maximizes expected reward
  - maximizes expected (reward-cost)

- given a ____ horizon
  - finite
  - infinite
  - indefinite

- assuming full observability

# Role of Discount Factor (γ)

- Keep the total reward/total cost finite
    - useful for infinite horizon problems

- Intuition (economics):
    - Money today is worth more than money tomorrow.

- Total reward: $r_1 + \gamma r_2 + \gamma^2 r_3 + \ldots$
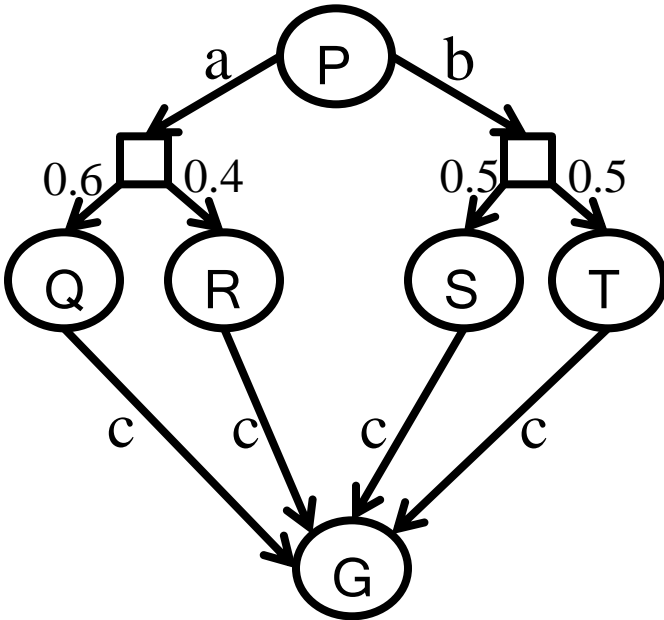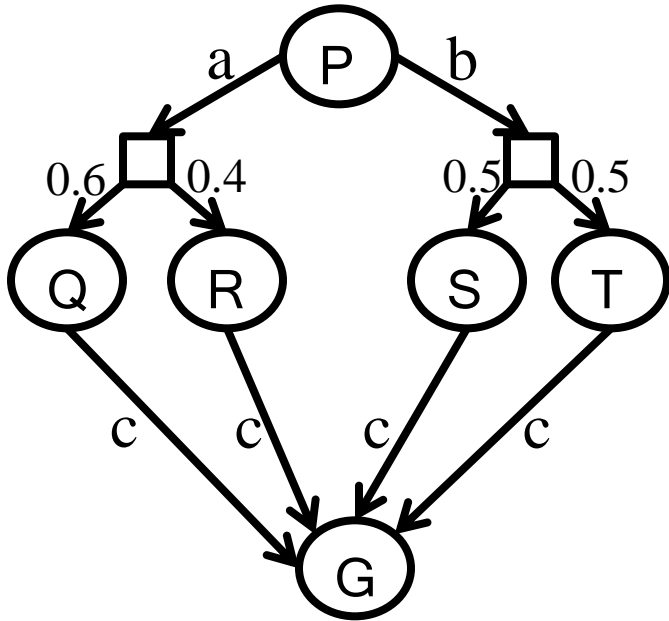- Total cost: $c_1 + \gamma c_2 + \gamma^2 c_3 + \ldots$

# Examples of MDPs

- Goal-directed, Indefinite Horizon, Cost Minimization MDP
  - $<\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{C}, \mathcal{G}, s_0>$
  - Most often studied in planning, graph theory communities

- Infinite Horizon, Discounted Reward Maximization MDP
  - $<\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma>$
  - Most often studied in machine learning, economics, operations research communities

- Oversubscription Planning: Non absorbing goals, Reward Max. MDP
  - $<\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{G}, \mathcal{R}, s_0>$
  - Relatively recent model

# Examples of MDPs

- Goal-directed, Indefinite Horizon, Cost Minimization MDP
  - $<\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{C}, \mathcal{G}, s_0>$
  - Most often studied in planning, graph theory communities

- Infinite Horizon, Discounted Reward Maximization MDP
  - $<\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma>$

    **most popular**
  - Most often studied in machine learning, economics, operations research communities

- Oversubscription Planning: Non absorbing goals, Reward Max. MDP
  - $<\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{G}, \mathcal{R}, s_0>$
  - Relatively recent model

# Acyclic vs. Cyclic MDPs



$$C(a) = 5, C(b) = 10, C(c) = 1$$

# Acyclic vs. Cyclic MDPs
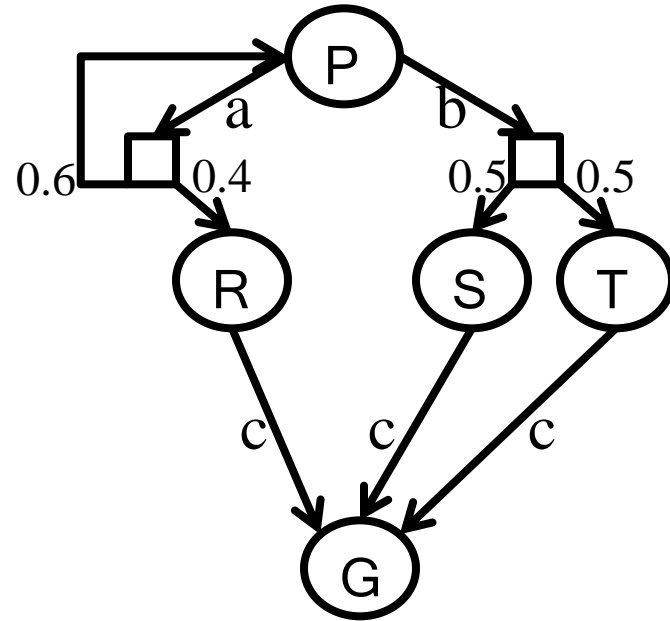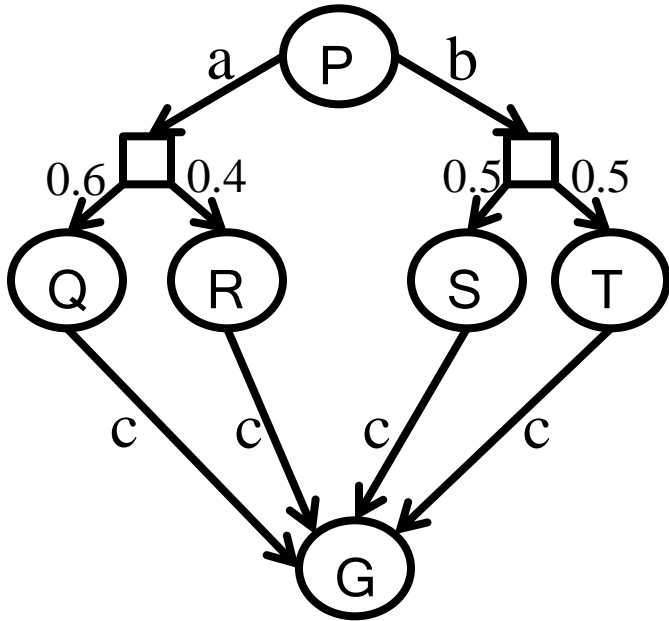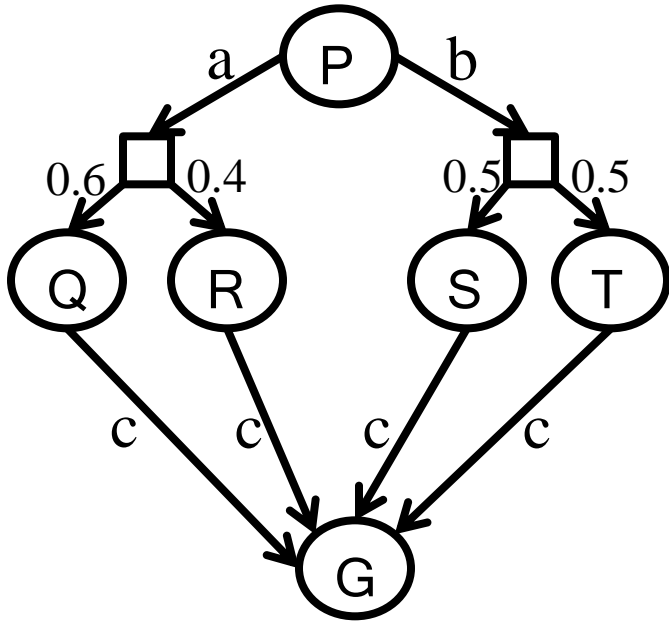


$C(a) = 5, C(b) = 10, C(c) = 1$

Expectimin works
- $V(Q/R/S/T) = 1$
- $V(P) = 6$ – action a

# Acyclic vs. Cyclic MDPs



C(a) = 5, C(b) = 10, C(c) =1

Expectimin works
- V(Q/R/S/T) = 1
- V(P) = 6 – action a

# Acyclic vs. Cyclic MDPs



C(a) = 5, C(b) = 10, C(c) =1

Expectimin doesn't work
•infinite loop
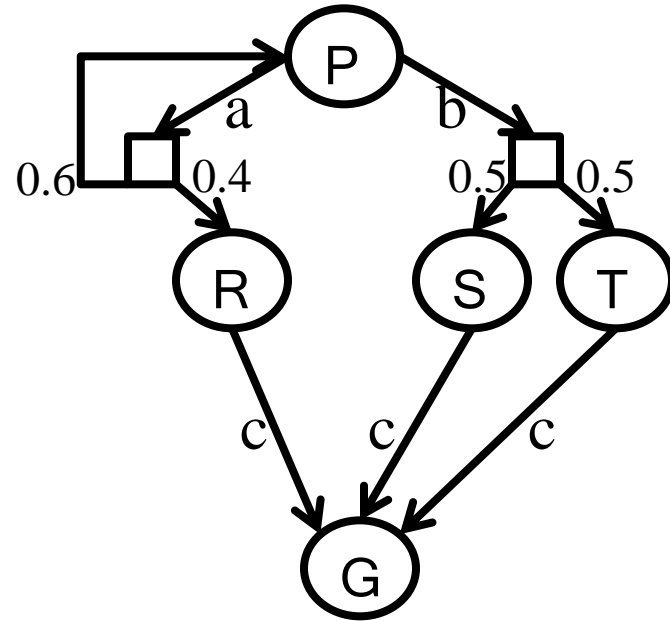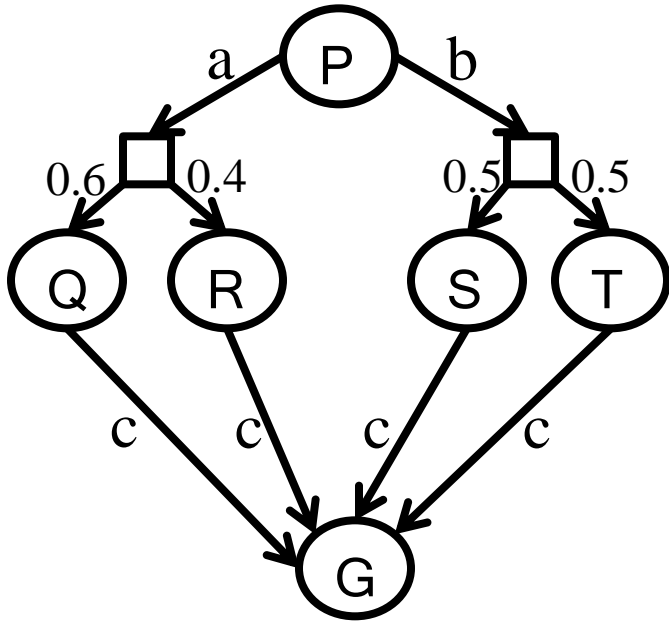
Expectimin works
• V(Q/R/S/T) = 1
• V(P) = 6 – action a

# Acyclic vs. Cyclic MDPs



C(a) = 5, C(b) = 10, C(c) = 1

Expectimin works
- V(Q/R/S/T) = 1
- V(P) = 6 – action a

Expectimin doesn't work
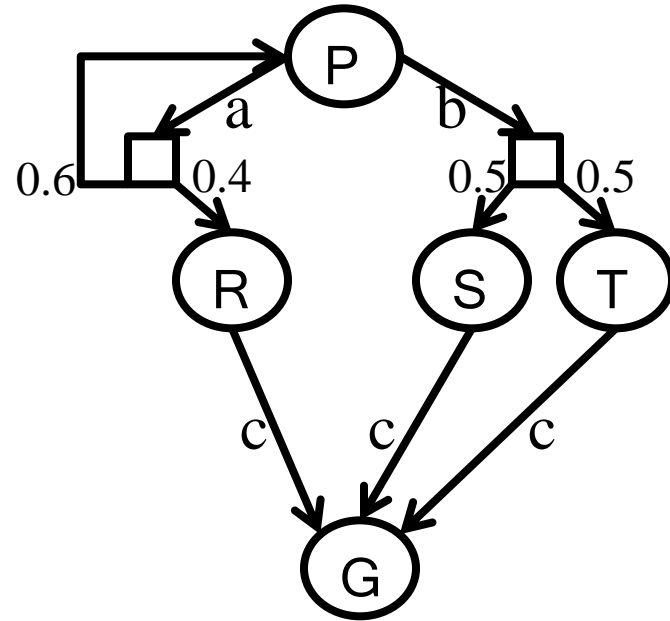  - infinite loop
- V(R/S/T) = 1
- Q(P,b) = 11
- Q(P,a) = ????

# Acyclic vs. Cyclic MDPs



C(a) = 5, C(b) = 10, C(c) =1

Expectimin works
- V(Q/R/S/T) = 1
- V(P) = 6 – action a

Expectimin doesn't work
  - infinite loop
- V(R/S/T) = 1
- Q(P,b) = 11
- Q(P,a) = ????
- suppose I decide to take a in P

# Acyclic vs. Cyclic MDPs



$C(a) = 5, C(b) = 10, C(c) = 1$

Expectimin works
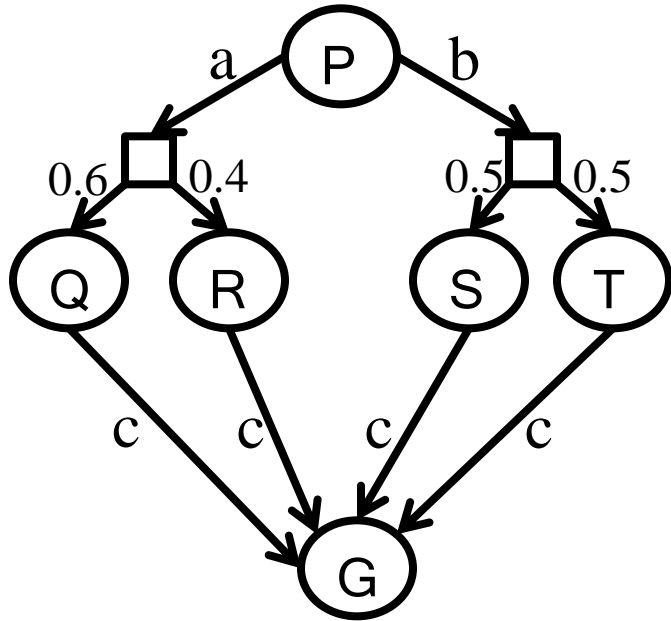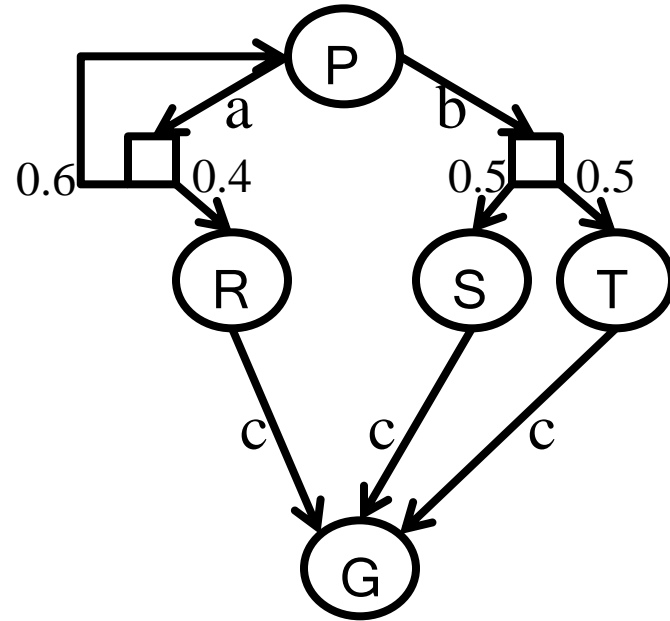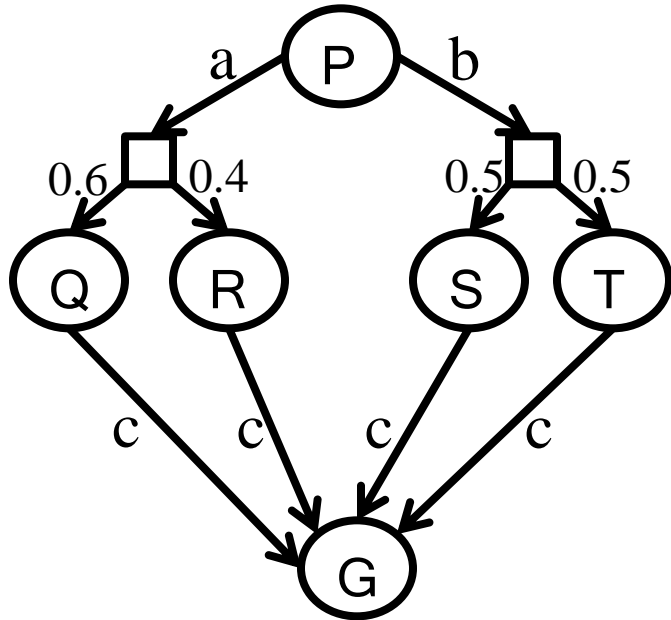- $V(Q/R/S/T) = 1$
- $V(P) = 6$ – action a

Expectimin doesn't work
   •infinite loop
- $V(R/S/T) = 1$
- $Q(P,b) = 11$
- $Q(P,a) = ????$
- suppose I decide to take a in P
- $Q(P,a) = 5 + 0.4*1 + 0.6Q(P,a)$
- ➔          $= 13.5$

# Policy Evaluation

- Given a policy $\pi$: compute $V^\pi$
  - $V^\pi$ : cost of reaching goal while following $\pi$

# Deterministic MDPs

- Policy Graph for $\pi$

    $\pi(s_0) = a_0;\ \pi(s_1) = a_1$



- $V^{\pi}(s_1) = 1$
- $V^{\pi}(s_0) = 6$

# Acyclic MDPs

- Policy Graph for $\pi$



- $V^\pi(s_1) = 1$
- $V^\pi(s_2) = 4$
- $V^\pi(s_0) = 0.6(5+1) + 0.4(2+4) = 6$

# General MDPs can be cyclic!



- $V^\pi(s_1) = 1$
- $V^\pi(s_2) = $ ?? (depends on $V^\pi(s_0)$)
- $V^\pi(s_0) = $ ?? (depends on $V^\pi(s_2)$)

# General SSPs can be cyclic!



- $V^\pi(g) = 0$
- $V^\pi(s_1) = 1 + V^\pi(s_g) = 1$
- $V^\pi(s_2) = 0.7(4 + V^\pi(s_g)) + 0.3(3 + V^\pi(s_0))$
- $V^\pi(s_0) = 0.6(5 + V^\pi(s_1)) + 0.4(2 + V^\pi(s_2))$

## Policy Evaluation (Approach 1)

- Solving the System of Linear Equations

$$V^\pi(s) \quad = \quad 0 \qquad \text{if } s \in \mathcal{G}$$

$$=$$

- $|S|$ variables.
- $O(|S|^3)$ running time

# Policy Evaluation (Approach 1)

- Solving the System of Linear Equations

$$V^\pi(s) \quad = \quad 0 \qquad \text{if } s \in \mathcal{G}$$

$$= \qquad\qquad\qquad [\mathcal{C}(s, \pi(s), s')$$

- $|S|$ variables.
- $\mathcal{O}(|S|^3)$ running time

# Policy Evaluation (Approach 1)

- Solving the System of Linear Equations

$$
\begin{aligned}
V^\pi(s) \quad &= \quad 0 \qquad \text{if } s \in \mathcal{G} \\[2ex]
&= \qquad\qquad\qquad [\mathcal{C}(s, \pi(s), s') + V^\pi(s')]
\end{aligned}
$$

- $|S|$ variables.
- $O(|S|^3)$ running time

# Policy Evaluation (Approach 1)

- Solving the System of Linear Equations

$$\begin{aligned} V^\pi(s) \quad &= \quad 0 \qquad \text{if } s \in \mathcal{G} \\ &= \quad \sum_{s' \in \mathcal{S}} \mathcal{T}(s, \pi(s), s') \left[ \mathcal{C}(s, \pi(s), s') + V^\pi(s') \right] \end{aligned}$$

- $|S|$ variables.
- $\mathcal{O}(|S|^3)$ running time

# Iterative Policy Evaluation

# Iterative Policy Evaluation



**1**

$S_1$

C=1

**0**

Pr=0.6
C=5

$a_1$

$S_g$

**4.4+0.4V$^\pi$(s$_2$)**

$S_0$

$a_0$

Pr=0.7
C=4

Pr=0.4
C=2

$S_2$

$a_2$

Pr=0.3
C=3

**3.7+0.3V$^\pi$(s$_0$)**

# Iterative Policy Evaluation



**1**

$S_1$

C=1

**0**

Pr=0.6
C=5

$a_1$

$S_g$

$a_0$

**4.4+0.4V$^\pi$(s$_2$)**
**0**

$S_0$

Pr=0.7
C=4

Pr=0.4
C=2

$S_2$

$a_2$

Pr=0.3
C=3

**3.7+0.3V$^\pi$(s$_0$)**

# Iterative Policy Evaluation

# Iterative Policy Evaluation

# Iterative Policy Evaluation



**4.4+0.4V$^\pi$(s$_2$)**
**0**
**5.88**

S$_0$

a$_0$

Pr=0.6
C=5

**1**

S$_1$

C=1

a$_1$

**0**

S$_g$

Pr=0.4
C=2

S$_2$

a$_2$

Pr=0.7
C=4

Pr=0.3
C=3

**3.7+0.3V$^\pi$(s$_0$)**
**3.7**
**5.464**

43

# Iterative Policy Evaluation



**1**

Pr=0.6
C=5

C=1

$a_1$

**0**

**4.4+0.4V$^\pi$(s$_2$)**
**0**
**5.88**
**6.5856**

$a_0$

$s_0$

$s_1$

$s_g$

Pr=0.4
C=2

$a_2$

Pr=0.7
C=4

$s_2$

Pr=0.3
C=3

**3.7+0.3V$^\pi$(s$_0$)**
**3.7**
**5.464**

# Iterative Policy Evaluation



**1**

$S_1$

C=1

Pr=0.6
C=5

$a_1$

**0**

$S_g$

$a_0$

**4.4+0.4V$^\pi$(s$_2$)**
**0**
**5.88**
**6.5856**

$S_0$

Pr=0.4
C=2

Pr=0.7
C=4

$S_2$

$a_2$

Pr=0.3
C=3

**3.7+0.3V$^\pi$(s$_0$)**
**3.7**
**5.464**
**5.67568**

# Iterative Policy Evaluation



**1**

**0**

$$C=1$$

$$a_1$$

**Pr=0.6**
**C=5**

$$a_0$$

$S_1$

$S_g$

**4.4+0.4V$^\pi$(s$_2$)**
**0**
**5.88**
**6.5856**
**6.670272**

$S_0$

**Pr=0.4**
**C=2**

$$a_2$$

**Pr=0.7**
**C=4**

$S_2$

**Pr=0.3**
**C=3**

**3.7+0.3V$^\pi$(s$_0$)**
**3.7**
**5.464**
**5.67568**

# Iterative Policy Evaluation



**1**

**0**

**4.4+0.4Vπ(s₂)**
**0**
**5.88**
**6.5856**
**6.670272**

$s_1$   C=1
$a_1$
Pr=0.6
C=5
$a_0$   $s_g$
$s_0$
Pr=0.4   Pr=0.7
C=2   C=4
$s_2$   $a_2$
Pr=0.3
C=3

**3.7+0.3Vπ(s₀)**
**3.7**
**5.464**
**5.67568**
**5.7010816**

# Iterative Policy Evaluation



**1**

**0**

**4.4+0.4V$^\pi$(s$_2$)**
**0**
**5.88**
**6.5856**
**6.670272**
**6.68043..**

S$_1$

C=1

a$_1$

S$_g$

Pr=0.6
C=5

a$_0$

S$_0$

Pr=0.4
C=2

Pr=0.7
C=4

S$_2$

a$_2$

Pr=0.3
C=3

**3.7+0.3V$^\pi$(s$_0$)**
**3.7**
**5.464**
**5.67568**
**5.7010816**

48

# Iterative Policy Evaluation



**4.4+0.4V$^\pi$(s$_2$)**
**0**
**5.88**
**6.5856**
**6.670272**
**6.68043..**

**1**

**0**

Pr=0.6
C=5

C=1

a$_1$

Pr=0.4
C=2

a$_2$

Pr=0.7
C=4

Pr=0.3
C=3

a$_0$

S$_0$

S$_1$

S$_2$

S$_g$

**3.7+0.3V$^\pi$(s$_0$)**
**3.7**
**5.464**
**5.67568**
**5.7010816**
**5.704129…**

49

# Policy Evaluation (Approach 2)

$$V^{\pi}(s) = \sum_{s' \in \mathcal{S}} \mathcal{T}(s, \pi(s), s') \left[ \mathcal{C}(s, \pi(s), s') + V^{\pi}(s') \right]$$

# Policy Evaluation (Approach 2)

$$V^\pi(s) = \sum_{s' \in \mathcal{S}} \mathcal{T}(s, \pi(s), s') \left[ \mathcal{C}(s, \pi(s), s') + V^\pi(s') \right]$$

*iterative refinement*

$$V_n^\pi(s) \leftarrow \sum_{s' \in \mathcal{S}} \mathcal{T}(s, \pi(s), s') \left[ \mathcal{C}(s, \pi(s), s') + V_{n-1}^\pi(s') \right]$$

# Policy Evaluation (Approach 2)

$$V^{\pi}(s) = \sum_{s' \in \mathcal{S}} \mathcal{T}(s, \pi(s), s') \left[ \mathcal{C}(s, \pi(s), s') + V^{\pi}(s') \right]$$

*iterative refinement*

$$V_n^{\pi}(s) \leftarrow \sum_{s' \in \mathcal{S}} \mathcal{T}(s, \pi(s), s') \left[ \mathcal{C}(s, \pi(s), s') + V_{n-1}^{\pi}(s') \right]$$

# Policy Evaluation (Approach 2)

$$V^\pi(s) = \sum_{s' \in \mathcal{S}} \mathcal{T}(s, \pi(s), s') \left[ \mathcal{C}(s, \pi(s), s') + V^\pi(s') \right]$$

*iterative refinement*

$$V^\pi_n(s) \leftarrow \sum_{s' \in \mathcal{S}} \mathcal{T}(s, \pi(s), s') \left[ \mathcal{C}(s, \pi(s), s') + V^\pi_{n-1}(s') \right]$$

# Iterative Policy Evaluation

```
 1   //Assumption: π is proper
 2   initialize V₀^π arbitrarily for each state
 3
 4
 5
 6
 7
 8
 9
10
11
```

# Iterative Policy Evaluation

1 $//Assumption: \pi \text{ is proper}$
2 initialize $V_0^\pi$ arbitrarily for each state
3
4
5
6
7
8
9
10
11

# Iterative Policy Evaluation

1  $//Assumption:\ \pi\ is\ proper$
2  initialize $V_0^\pi$ arbitrarily for each state
3  $n \leftarrow 0$
4  **repeat**
5      $n \leftarrow n + 1$
6      **foreach** $s \in \mathcal{S}$ **do**
7          compute $V_n^\pi(s) \leftarrow \sum_{s' \in \mathcal{S}} \mathcal{T}(s, \pi(s), s') \left[ \mathcal{C}(s, \pi(s), s') + V_{n-1}^\pi(s') \right]$
8
9      **end**
10
11

# Iterative Policy Evaluation

1 *//Assumption: $\pi$ is proper*
2 initialize $V_0^\pi$ arbitrarily for each state
3 $n \leftarrow 0$
4 **repeat**
5     $n \leftarrow n + 1$
6     **foreach** $s \in \mathcal{S}$ **do**
7         compute $V_n^\pi(s) \leftarrow \sum_{s' \in \mathcal{S}} \mathcal{T}(s, \pi(s), s') \left[ \mathcal{C}(s, \pi(s), s') + V_{n-1}^\pi(s') \right]$
8
9     **end**
10
11

*iteration n*

57

# Iterative Policy Evaluation

1 $//Assumption: \pi \ is \ proper$
2 initialize $V_0^\pi$ arbitrarily for each state
3 $n \leftarrow 0$
4 **repeat**
5 $\quad$ $n \leftarrow n + 1$
6 $\quad$ **foreach** $s \in \mathcal{S}$ **do**
7 $\quad\quad$ compute $V_n^\pi(s) \leftarrow \sum_{s' \in \mathcal{S}} \mathcal{T}(s, \pi(s), s') \left[ \mathcal{C}(s, \pi(s), s') + V_{n-1}^\pi(s') \right]$
8 $\quad\quad$ compute $\text{residual}_n(s) \leftarrow |V_n^\pi(s) - V_{n-1}^\pi(s)|$
9 $\quad$ **end**
10
11

# Iterative Policy Evaluation

1 $//Assumption: \pi$ is proper
2 initialize $V_0^\pi$ arbitrarily for each state
3 $n \leftarrow 0$
4 **repeat**
5      $n \leftarrow n + 1$
6      **foreach** $s \in \mathcal{S}$ **do**
7          compute $V_n^\pi(s) \leftarrow \sum_{s' \in \mathcal{S}} \mathcal{T}(s, \pi(s), s') \left[ \mathcal{C}(s, \pi(s), s') + V_{n-1}^\pi(s') \right]$
8          compute $\text{residual}_n(s) \leftarrow |V_n^\pi(s) - V_{n-1}^\pi(s)|$
9      **end**
10 **until** $\max_{s \in \mathcal{S}} \text{residual}_n(s) < \epsilon$;
11 return $V_n^\pi$

# Iterative Policy Evaluation

1 //*Assumption: $\pi$ is proper*
2 initialize $V_0^\pi$ arbitrarily for each state
3 $n \leftarrow 0$
4 **repeat**
5      $n \leftarrow n + 1$
6      **foreach** $s \in \mathcal{S}$ **do**
7          compute $V_n^\pi(s) \leftarrow \sum_{s' \in \mathcal{S}} \mathcal{T}(s, \pi(s), s') \left[ \mathcal{C}(s, \pi(s), s') + V_{n-1}^\pi(s') \right]$
8          compute $\text{residual}_n(s) \leftarrow |V_n^\pi(s) - V_{n-1}^\pi(s)|$
9      **end**
10 **until** $\max_{s \in \mathcal{S}} \text{residual}_n(s) < \epsilon;$    *$\epsilon$-consistency*
11 return $V_n^\pi$

60

# Iterative Policy Evaluation

1  $//Assumption: \pi\ is\ proper$
2  initialize $V_0^\pi$ arbitrarily for each state
3  $n \leftarrow 0$
4  **repeat**
5      $n \leftarrow n + 1$
6      **foreach** $s \in \mathcal{S}$ **do**
7          compute $V_n^\pi(s) \leftarrow \sum_{s' \in \mathcal{S}} \mathcal{T}(s, \pi(s), s') \left[ \mathcal{C}(s, \pi(s), s') + V_{n-1}^\pi(s') \right]$
8          compute $\text{residual}_n(s) \leftarrow |V_n^\pi(s) - V_{n-1}^\pi(s)|$
9      **end**
10 **until** $\max_{s \in \mathcal{S}} \text{residual}_n(s) < \epsilon;$
11 return $V_n^\pi$

*$\epsilon$-consistency*

*termination condition*

61

# Policy Evaluation → Value Iteration (Bellman Equations for MDP$_1$)

- $<\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{C}, \mathcal{G}, s_0>$

- Define V*(s) {optimal cost} as the minimum expected cost to reach a goal from this state.

- V* should satisfy the following equation:

$$
\begin{aligned}
V^*(s) \;=\;& 0 \qquad \text{if } s \in \mathcal{G} \\
=\;&
\end{aligned}
$$

# Policy Evaluation → Value Iteration (Bellman Equations for MDP₁)

- $\langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{C}, \mathcal{G}, s_0 \rangle$

- Define V*(s) {optimal cost} as the minimum expected cost to reach a goal from this state.

- V* should satisfy the following equation:

$$
\begin{aligned}
V^*(s) &= 0 &&\text{if } s \in \mathcal{G} \\
&= && [\mathcal{C}(s, a, s')
\end{aligned}
$$

# Policy Evaluation ➔ Value Iteration (Bellman Equations for MDP$_1$)

- $<\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{C}, \mathcal{G}, s_0>$

- Define V*(s) {optimal cost} as the minimum expected cost to reach a goal from this state.

- V* should satisfy the following equation:

$$V^*(s) \quad = \quad 0 \qquad \text{if } s \in \mathcal{G}$$

$$= \qquad\qquad\qquad\qquad [\mathcal{C}(s, a, s') + V^*(s')]$$

# Policy Evaluation → Value Iteration (Bellman Equations for MDP$_1$)

- $<\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{C}, \mathcal{G}, s_0>$

- Define V*(s) {optimal cost} as the minimum expected cost to reach a goal from this state.

- V* should satisfy the following equation:

$$
\begin{aligned}
V^*(s) \quad &= \quad 0 \qquad \text{if } s \in \mathcal{G} \\
&= \quad \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V^*(s') \right]
\end{aligned}
$$

# Policy Evaluation → Value Iteration (Bellman Equations for MDP$_1$)

- $<\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{C}, \mathcal{G}, s_0>$

- Define V*(s) {optimal cost} as the minimum expected cost to reach a goal from this state.

- V* should satisfy the following equation:

$$
\begin{aligned}
V^*(s) &= 0 && \text{if } s \in \mathcal{G} \\
&= \min_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V^*(s') \right]
\end{aligned}
$$

# Policy Evaluation ➔ Value Iteration (Bellman Equations for MDP$_1$)

- $<\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{C}, \mathcal{G}, s_0>$

- Define V*(s) {optimal cost} as the minimum expected cost to reach a goal from this state.

- V* should satisfy the following equation:

$$
\begin{aligned}
V^*(s) &= 0 && \text{if } s \in \mathcal{G} \\
&= \min_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V^*(s') \right]
\end{aligned}
$$

## Policy Evaluation → Value Iteration (Bellman Equations for MDP$_1$)

- $<\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{C}, \mathcal{G}, s_0>$

- Define V*(s) {optimal cost} as the minimum expected cost to reach a goal from this state.

- V* should satisfy the following equation:

$$
\begin{aligned}
V^*(s) &= 0 && \text{if } s \in \mathcal{G} \\
&= \min_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V^*(s') \right]
\end{aligned}
$$

**Q*(s,a)**

**V*(s) = min$_a$ Q*(s,a)**

# Bellman Equations for MDP$_2$

- $<\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, s_0, \gamma>$

- Define **V*(s)** {optimal **value**} as the **maximum** expected **discounted reward** from this state.

- V* should satisfy the following equation:

$$V^*(s) \;=\; \max_{a \in Ap(s)} \sum_{s' \in \mathcal{S}} Pr(s'|s,a) \left[ \mathcal{R}(s,a,s') + \gamma V^*(s') \right]$$
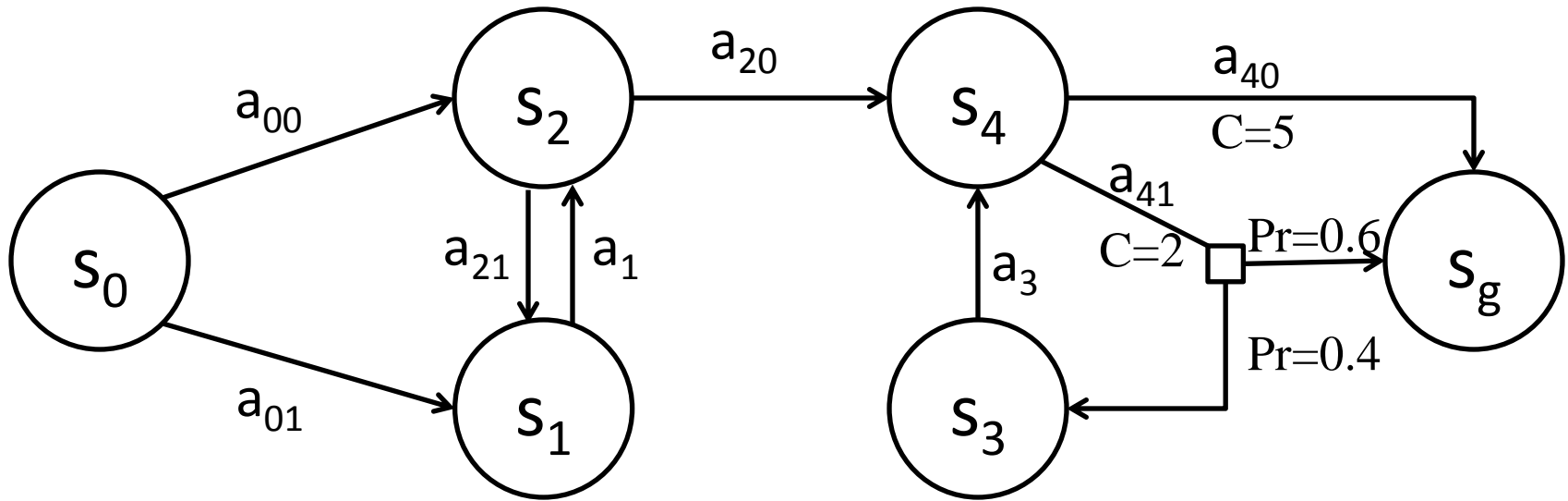
# Fixed Point Computation in VI

$$V^*(s) = \min_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V^*(s') \right]$$
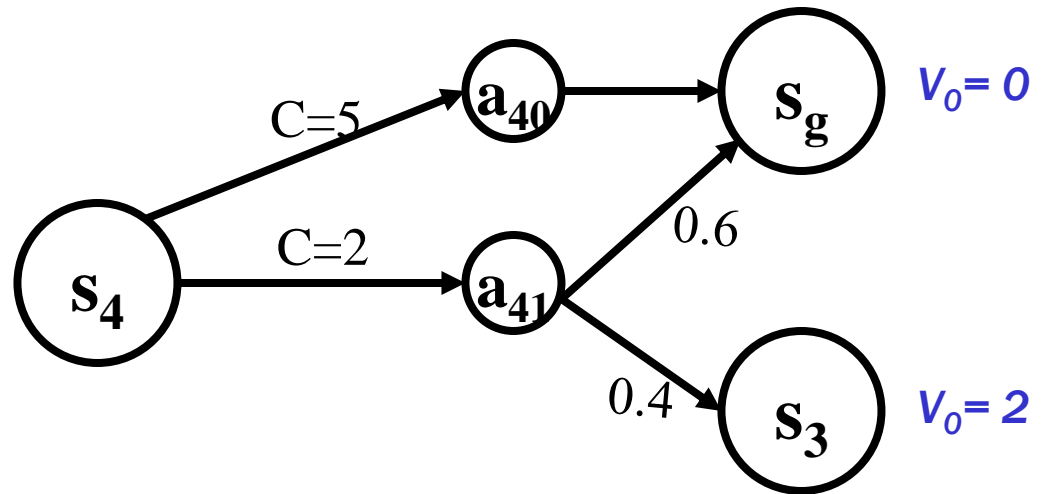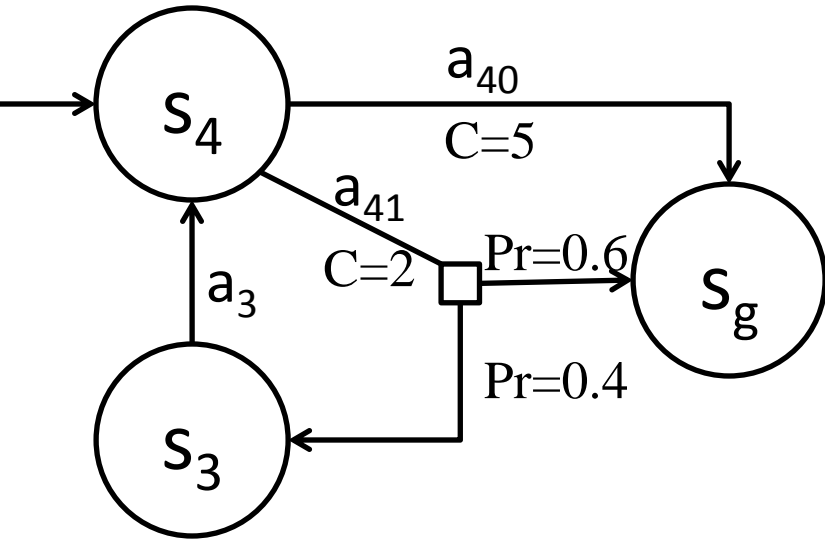
*iterative refinement*

$$V_n(s) \leftarrow \min_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V_{n-1}(s') \right]$$

# Fixed Point Computation in VI

$$V^*(s) = \min_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V^*(s') \right]$$

*iterative refinement*

$$V_n(s) \leftarrow \min_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V_{n-1}(s') \right]$$
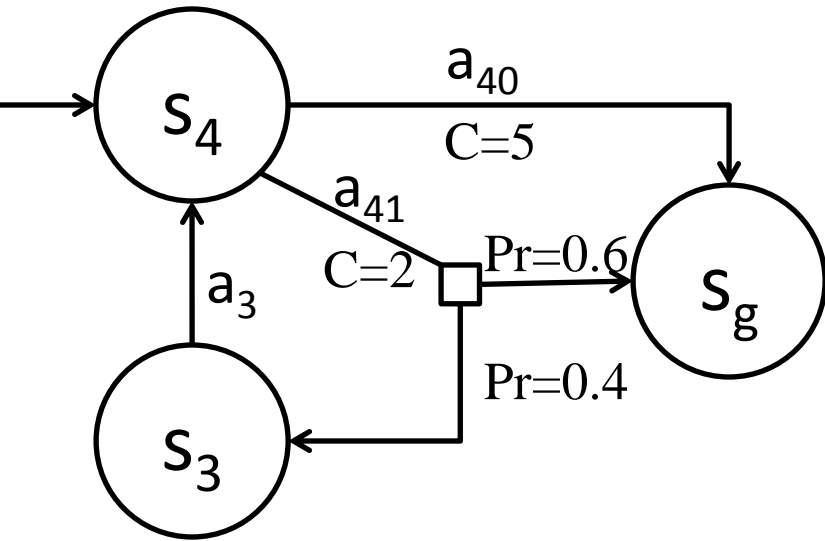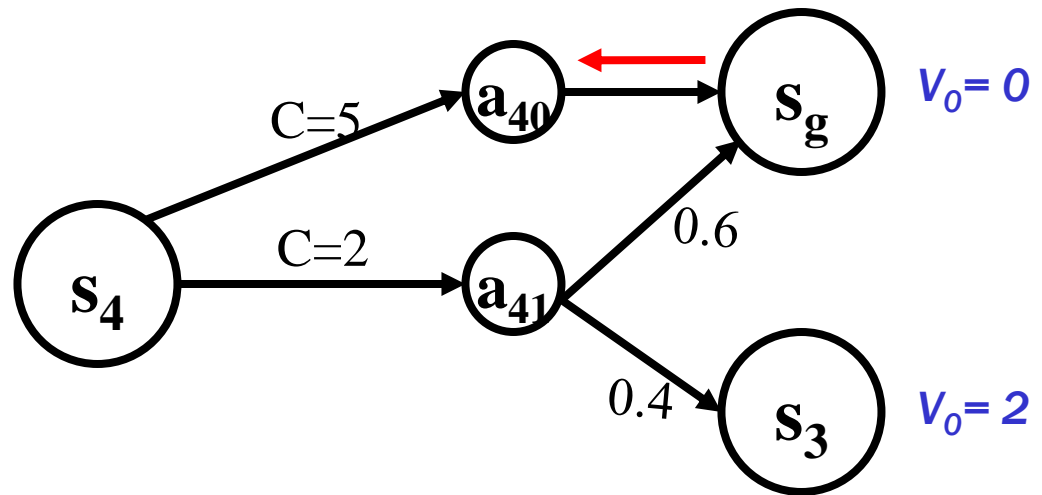
# Fixed Point Computation in VI

$$V^*(s) = \min_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V^*(s') \right]$$

*iterative refinement*

$$V_n(s) \leftarrow \min_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V_{n-1}(s') \right]$$

# Fixed Point Computation in VI

$$V^*(s) = \min_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V^*(s') \right]$$

*iterative refinement*

$$V_n(s) \leftarrow \min_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V_{n-1}(s') \right]$$
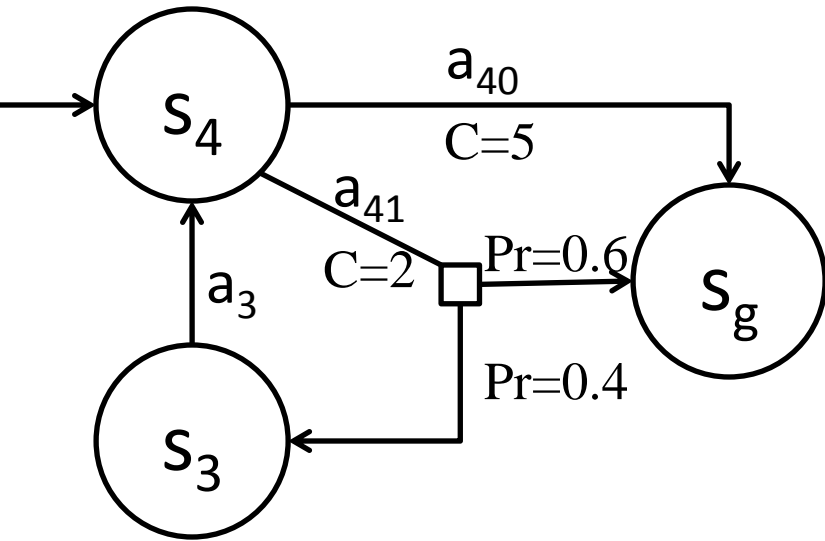
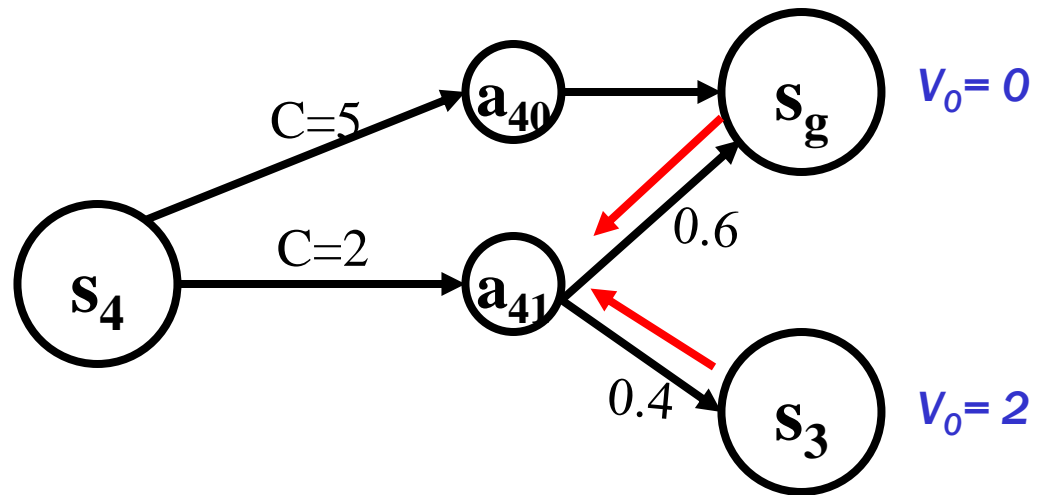*non-linear*
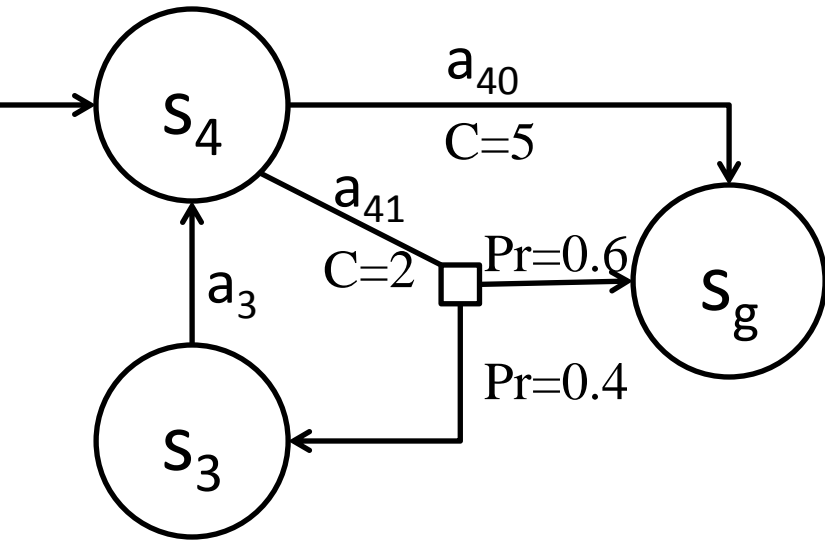
# Example

# Bellman Backup
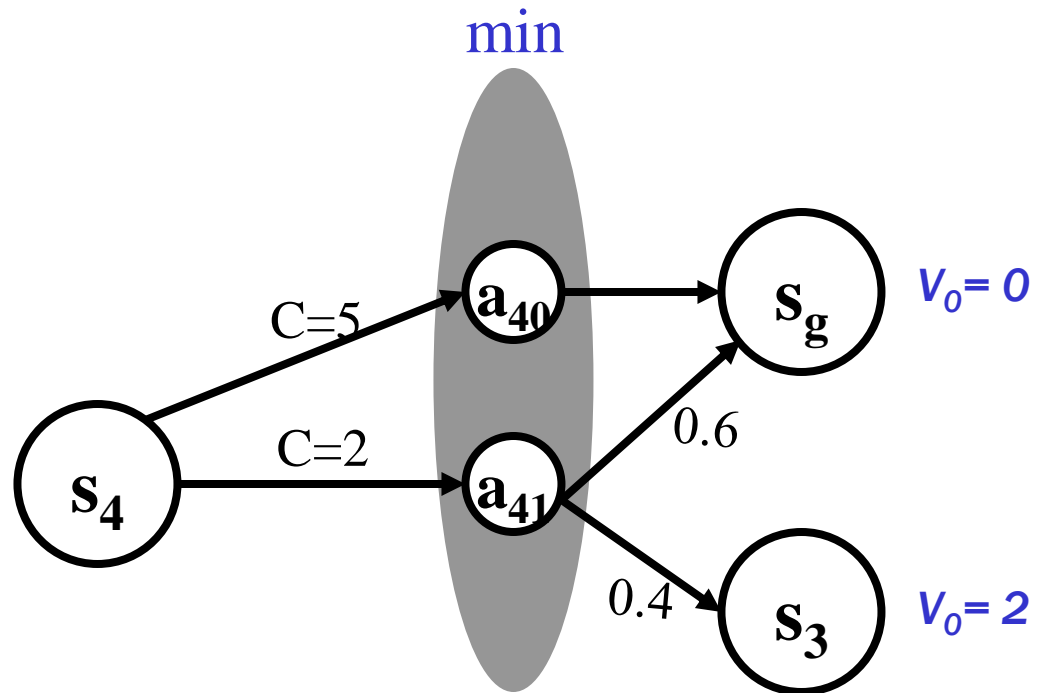
# Bellman Backup



$Q_1(s_4, a_{40}) = 5 + 0$

# Bellman Backup



$Q_1(s_4, a_{40}) = 5 + 0$

$Q_1(s_4, a_{41}) = 2 + 0.6 \times 0$
$$+ 0.4 \times 2$$
$$= 2.8$$

# Bellman Backup



$Q_1(s_4, a_{40}) = 5 + 0$
$Q_1(s_4, a_{41}) = 2 + 0.6 \times 0$
$+ 0.4 \times 2$
$= 2.8$

# Bellman Backup



$Q_1(s_4, a_{40}) = 5 + 0$

$Q_1(s_4, a_{41}) = 2 + 0.6 \times 0$
$+ 0.4 \times 2$
$= 2.8$

# Bellman Backup



$Q_1(s_4, a_{40}) = 5 + 0$

$Q_1(s_4, a_{41}) = 2 + 0.6 \times 0$
$+ 0.4 \times 2$
$= 2.8$

$a_3$

$s_3$

min

$a_{greedy} = a_{41}$

$V_1 = 2.8$

$C=5$

$C=2$

$a_{40}$

$a_{41}$

$s_g$   $V_0 = 0$

$s_4$

$0.6$

$0.4$

$s_3$   $V_0 = 2$

# Value Iteration [Bellman 57]

1   initialize $V_0$ arbitrarily for each state

2   $n \leftarrow 0$

3   **repeat**

4      $n \leftarrow n + 1$

5      **foreach** $s \in \mathcal{S}$ **do**

6         compute $V_n(s)$ using Bellman backup at $s$

7         compute $\text{residual}_n(s) = |V_n(s) - V_{n-1}(s)|$

8      **end**

9   **until** $\max_{s \in \mathcal{S}} \text{residual}_n(s) < \epsilon$;

10   return greedy policy: $\pi^{V_n}(s) = \text{argmin}_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V_n(s') \right]$

# Value Iteration [Bellman 57]

**No restriction on initial value function**

1   initialize $V_0$ arbitrarily for each state

2   $n \leftarrow 0$

3   **repeat**

4      $n \leftarrow n + 1$

5      **foreach** $s \in \mathcal{S}$ **do**

6         compute $V_n(s)$ using Bellman backup at $s$

7         compute $\text{residual}_n(s) = |V_n(s) - V_{n-1}(s)|$

8      **end**

9   **until** $\max_{s \in \mathcal{S}} \text{residual}_n(s) < \epsilon$;

10  return greedy policy: $\pi^{V_n}(s) = \text{argmin}_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s')\left[\mathcal{C}(s, a, s') + V_n(s')\right]$

# Value Iteration [Bellman 57]

1 initialize $V_0$ arbitrarily for each state
2 $n \leftarrow 0$
3 **repeat**
4 $\quad n \leftarrow n+1$
5 $\quad$ **foreach** $s \in \mathcal{S}$ **do**
6 $\qquad$ compute $V_n(s)$ using Bellman backup at $s$
7 $\qquad$ compute $\text{residual}_n(s) = |V_n(s) - V_{n-1}(s)|$
8 $\quad$ **end**
9 **until** $\max_{s \in \mathcal{S}} \text{residual}_n(s) < \epsilon$;
10 return greedy policy: $\pi^{V_n}(s) = \text{argmin}_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V_n(s') \right]$

*iteration n*

# Value Iteration [Bellman 57]

1   initialize $V_0$ arbitrarily for each state

2   $n \leftarrow 0$

3   **repeat**

4      $n \leftarrow n + 1$

5      **foreach** $s \in \mathcal{S}$ **do**

6         compute $V_n(s)$ using Bellman backup at $s$

7         compute $\text{residual}_n(s) = |V_n(s) - V_{n-1}(s)|$

8      **end**

9   **until** $\max_{s \in \mathcal{S}} \text{residual}_n(s) < \epsilon$;   &larr; *ε-consistency*

10   return greedy policy: $\pi^{V_n}(s) = \text{argmin}_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') \left[ \mathcal{C}(s, a, s') + V_n(s') \right]$
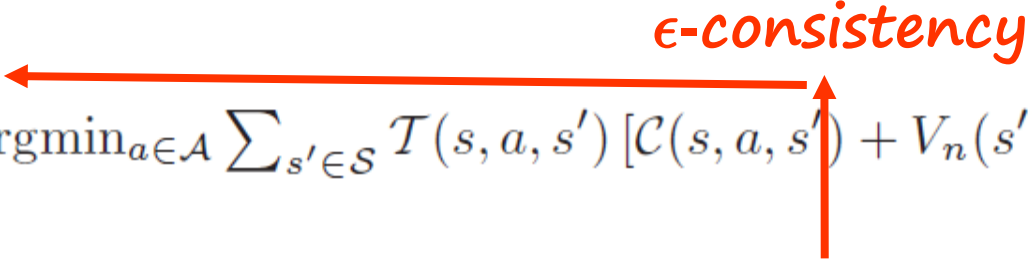
# Value Iteration [Bellman 57]

1   initialize $V_0$ arbitrarily for each state

2   $n \leftarrow 0$

3   **repeat**

4     $n \leftarrow n + 1$

5     **foreach** $s \in \mathcal{S}$ **do**

6       compute $V_n(s)$ using Bellman backup at $s$

7       compute $\mathrm{residual}_n(s) = |V_n(s) - V_{n-1}(s)|$

8     **end**

9   **until** $\max_{s \in \mathcal{S}} \mathrm{residual}_n(s) < \epsilon$;

10 return greedy policy: $\pi^{V_n}(s) = \mathrm{argmin}_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s')\left[\mathcal{C}(s, a, s') + V_n(s')\right]$
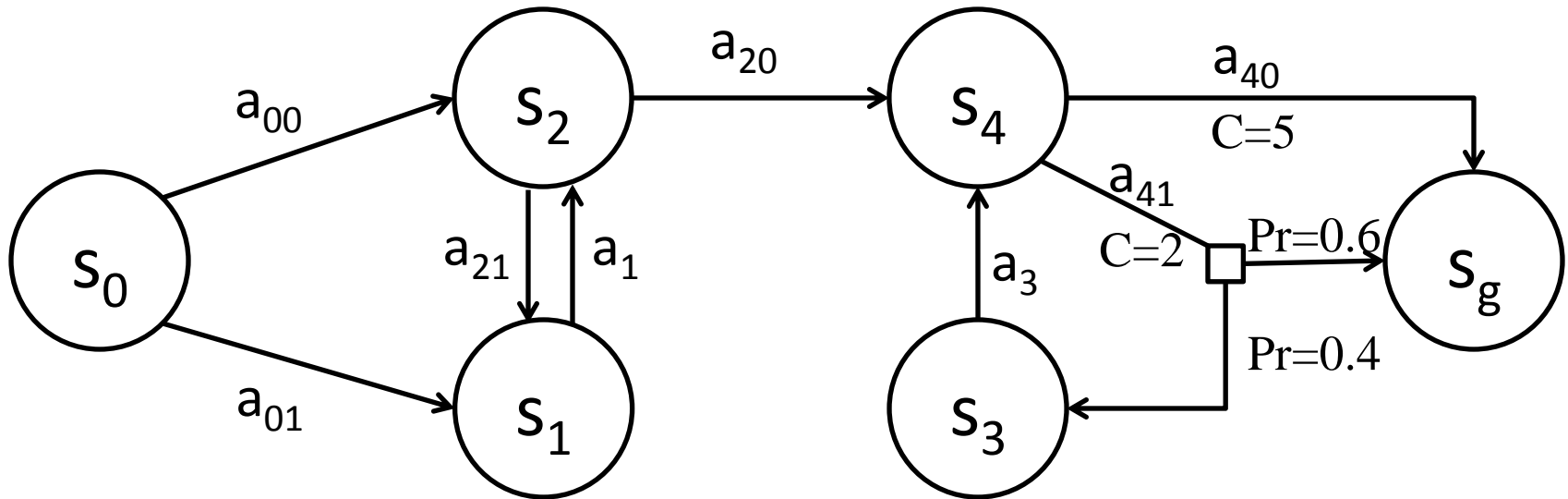
*$\epsilon$-consistency*

*termination condition*

85

# Example

(all actions cost 1 unless otherwise stated)



| n | $V_n(s_0)$ | $V_n(s_1)$ | $V_n(s_2)$ | $V_n(s_3)$ | $V_n(s_4)$ |
|---|---|---|---|---|---|
| 0 | 3 | 3 | 2 | 2 | 1 |
| 1 | 3 | 3 | 2 | 2 | 2.8 |
| 2 | 3 | 3 | 3.8 | 3.8 | 2.8 |
| 3 | 4 | 4.8 | 3.8 | 3.8 | 3.52 |
| 4 | 4.8 | 4.8 | 4.52 | 4.52 | 3.52 |
| 5 | 5.52 | 5.52 | 4.52 | 4.52 | 3.808 |
| 20 | 5.99921 | 5.99921 | 4.99969 | 4.99969 | 3.99969 |

# Comments

- Decision-theoretic Algorithm
- Dynamic Programming
- Fixed Point Computation
- Probabilistic version of Bellman-Ford Algorithm
  - for shortest path computation
  - $MDP_1$ : Stochastic Shortest Path Problem

- Time Complexity
  - one iteration: $O(|\mathcal{S}|^2|\mathcal{A}|)$
  - number of iterations: $poly(|\mathcal{S}|, |\mathcal{A}|, 1/(1-\gamma))$
- Space Complexity: $O(|\mathcal{S}|)$

# Monotonicity

For all n>k

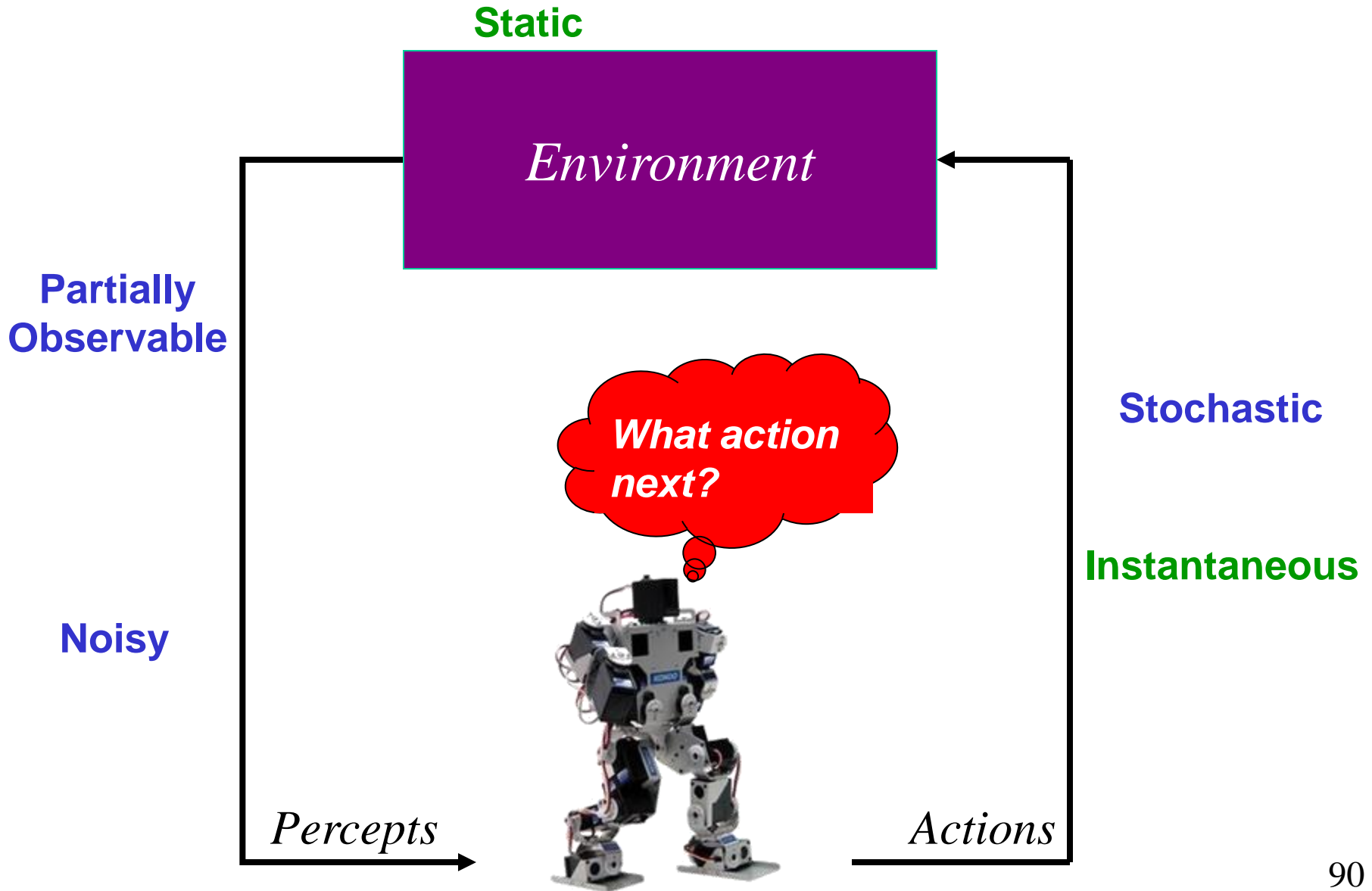$V_k \leq_p V^* \Rightarrow V_n \leq_p V^*$ ($V_n$ monotonic from below)

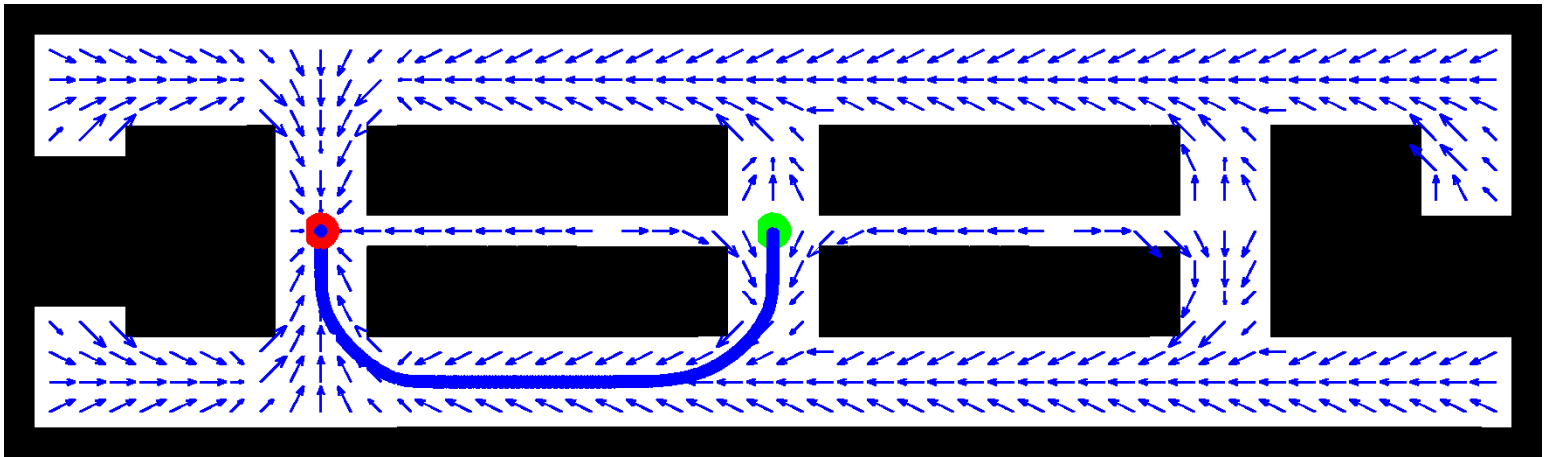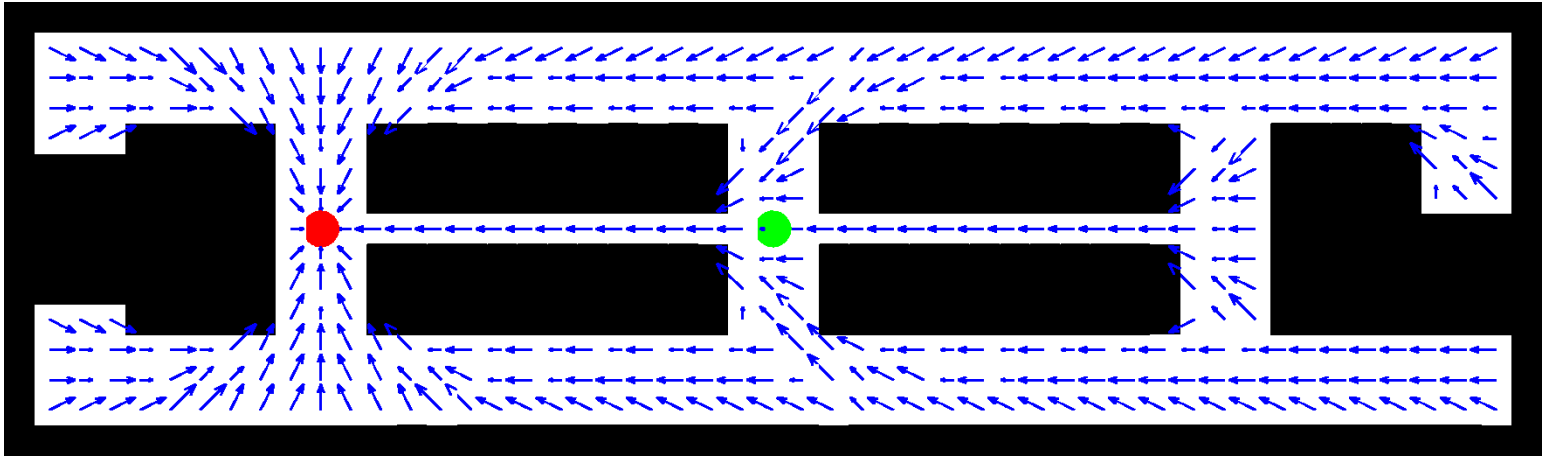$V_k \geq_p V^* \Rightarrow V_n \geq_p V^*$ ($V_n$ monotonic from above)

# Extensions

- Heuristic Search + Dynamic Programming
  - AO*, LAO*, RTDP, …

- Factored MDPs
  - add planning graph style heuristics
  - use goal regression to generalize better

- Hierarchical MDPs
  - hierarchy of sub-tasks, actions to scale better

- Reinforcement Learning
  - learning the probability and rewards
  - acting while learning – connections to psychology

- Partially Observable Markov Decision Processes
  - noisy sensors; partially observable environment
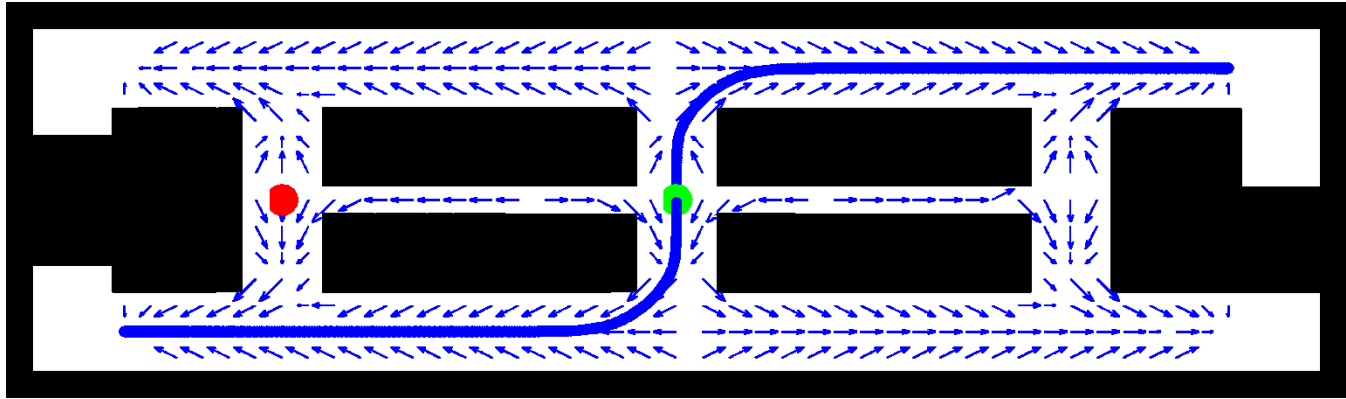  - popular in robotics

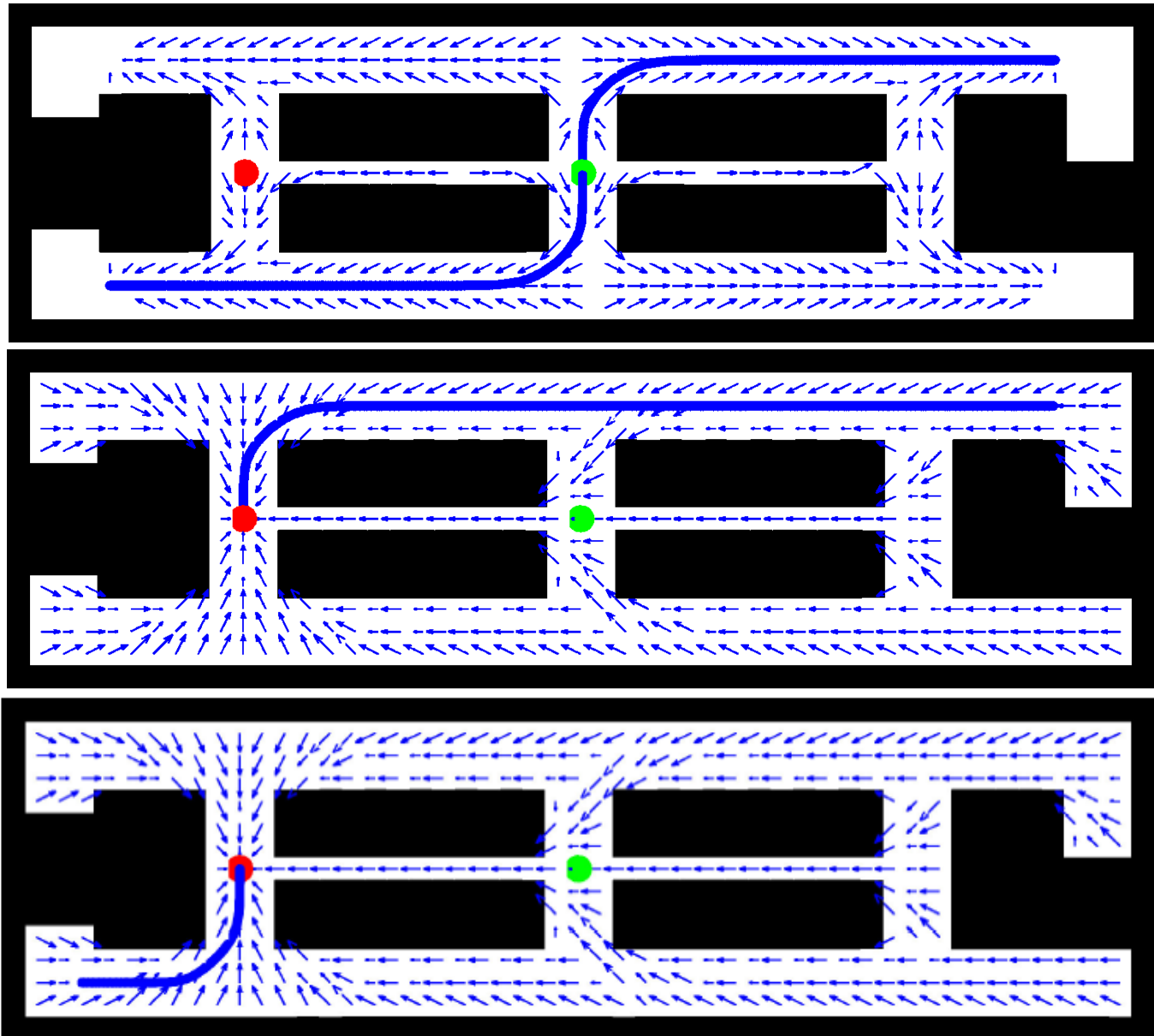# Partially Observable MDPs



Static

Environment

Partially
Observable

Stochastic

*What action next?*

Instantaneous

Noisy

*Percepts*

*Actions*

# Stochastic, Fully Observable

# Stochastic, Partially Observable

# Stochastic, Partially Observable

# POMDPs

- In POMDPs we apply the very same idea as in MDPs.

- Since the state is not observable,
    the agent has to make its decisions based on the belief state which is a posterior distribution over states.

- Let $b$ be the belief of the agent about the current state

- POMDPs compute a value function over belief space:

$$V_T(b) \;=\; \max_a \left[ r(b,a) + \gamma \int V_{T-1}(b')p(b' \mid b, a) \; db' \right]$$

# POMDPs

- Each belief is a probability distribution,

  - value fn is a function of an entire probability distribution.

- Problematic, since probability distributions are continuous.

- Also, we have to deal with huge complexity of belief spaces.

- For finite worlds with finite state, action, and observation spaces and finite horizons,

  - we can represent the value functions by piecewise linear functions.

# Applications

- Robotic control
  - helicopter maneuvering, autonomous vehicles
  - Mars rover - path planning, oversubscription planning
  - elevator planning
- Game playing - backgammon, tetris, checkers
- Neuroscience
- Computational Finance, Sequential Auctions
- Assisting elderly in simple tasks
- Spoken dialog management
- Communication Networks – switching, routing, flow control
- War planning, evacuation planning