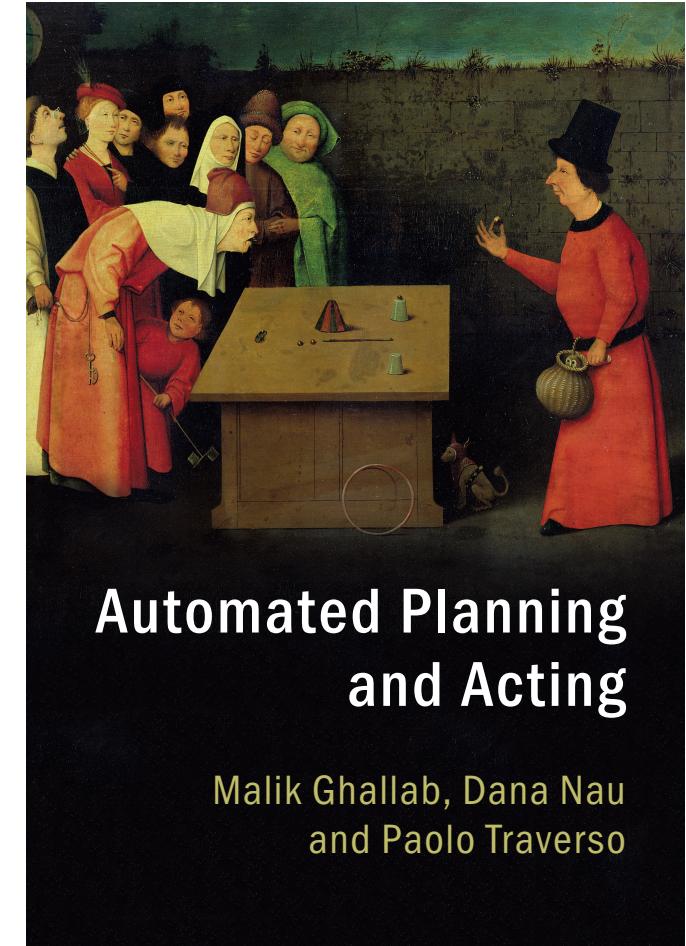


Chapter 6

Deliberation with Probabilistic Domain Models

Dana S. Nau
University of Maryland



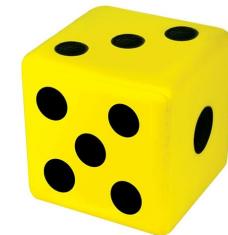
Automated Planning
and Acting

Malik Ghallab, Dana Nau
and Paolo Traverso

<http://www.laas.fr/planning>

Motivation

- Situations where actions have multiple possible outcomes and each outcome has a probability
- Several possible action representations
 - ▶ Bayes nets, probabilistic actions, ...
- Book doesn't commit to any representation
 - ▶ Mainly concentrates on the underlying semantics



$\text{roll-die}(d)$

pre: $\text{holding}(d) = \text{true}$
eff:

1/6: $\text{top}(d) \leftarrow 1$
1/6: $\text{top}(d) \leftarrow 2$
1/6: $\text{top}(d) \leftarrow 3$
1/6: $\text{top}(d) \leftarrow 4$
1/6: $\text{top}(d) \leftarrow 5$
1/6: $\text{top}(d) \leftarrow 6$

Probabilistic planning domain

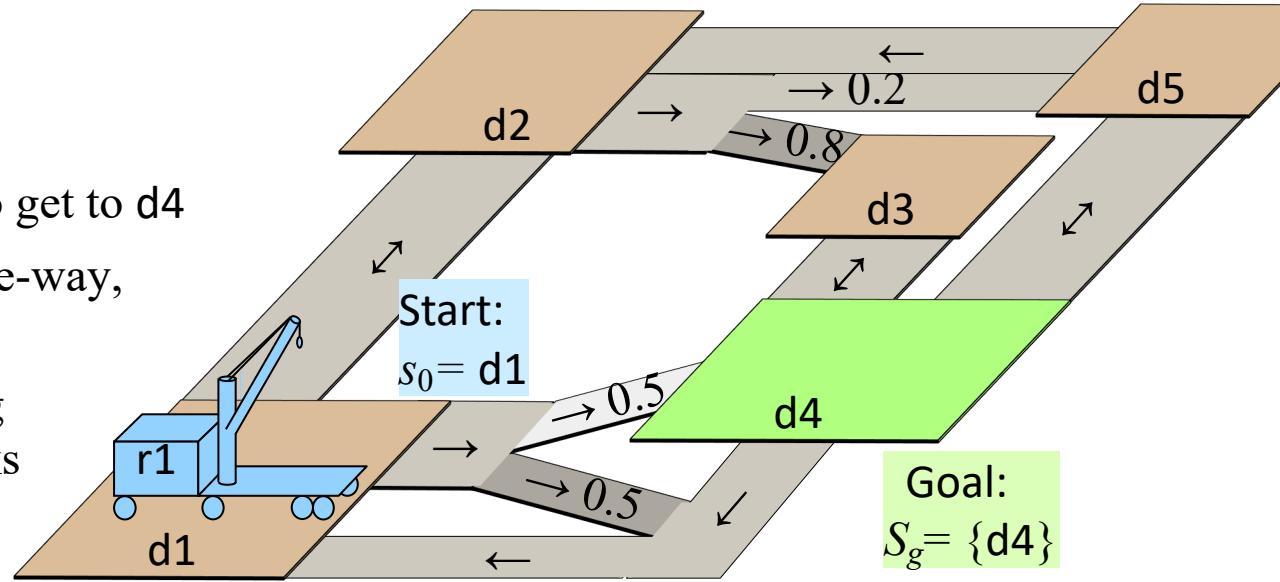
Definitions

$$\Sigma = (S, A, \gamma, \Pr, \text{cost})$$

- $S = \{\text{states}\}$
- $A = \{\text{actions}\}$
- $\gamma : S \times A \rightarrow 2^S$
- $\Pr(s' | s, a) = \text{probability of going to state } s' \text{ if we're in } s \text{ and apply } a$
 - ▶ $\Pr(s' | s, a) \neq 0 \text{ iff } s' \in \gamma(s, a)$
- cost: $S \times A \rightarrow \mathbb{R}_{\geq 0}$
 - ▶ $\text{cost}(s, a) = \text{cost of action } a \text{ in state } s$
 - ▶ may omit, default is $\text{cost}(s, a) = 1$
- Applicable(s) = $\{a \mid \gamma(s, a) \neq \emptyset\}$

Example

- Start at d_1 , want to get to d_4
- Some roads are one-way, some are two-way
- Unreliable steering when the road forks
 - ▶ may take the wrong fork
- Simplified state and action names:
 - ▶ write $\{\text{loc}(r1)=d_2\}$ as d_2
 - ▶ write $\text{move}(r1, d_2, d_3)$ as m_{23}
- $\gamma(d_1, m_{12}) = \{d_2\}$
 - ▶ $\Pr(d_2 | d_1, m_{12}) = 1$
- $m_{21}, m_{34}, m_{41}, m_{43}, m_{45}, m_{52}, m_{54}$:
 - ▶ like m_{12}



- $\gamma(d_1, m_{14}) = \{d_1, d_4\}$
 - ▶ $\Pr(d_4 | d_1, m_{14}) = 0.5$
 - ▶ $\Pr(d_1 | d_1, m_{14}) = 0.5$
- $\gamma(d_2, m_{23}) = \{d_3, d_5\}$
 - ▶ $\Pr(d_3 | d_2, m_{23}) = 0.8$
 - ▶ $\Pr(d_5 | d_2, m_{23}) = 0.2$
- No m_{11} , no m_{25}

Probabilistic planning domain

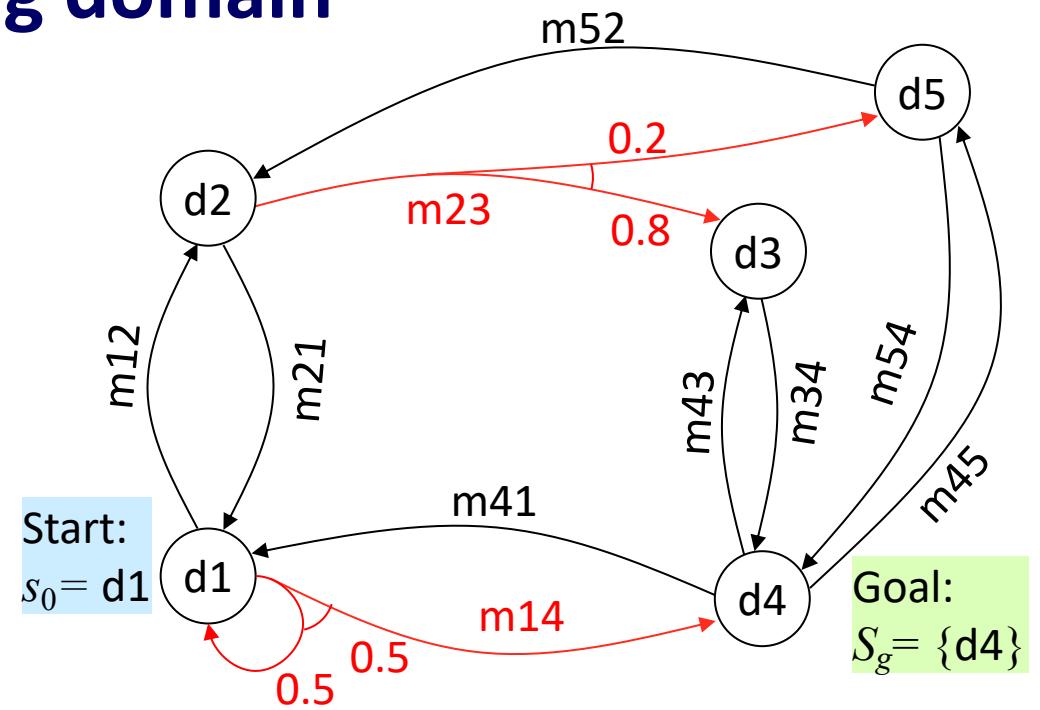
Definitions

$$\Sigma = (S, A, \gamma, \Pr, \text{cost})$$

- $S = \{\text{states}\}$
- $A = \{\text{actions}\}$
- $\gamma : S \times A \rightarrow 2^S$
- $\Pr(s' | s, a) = \text{probability of going to state } s' \text{ if we're in } s \text{ and apply } a$
 - ▶ $\Pr(s' | s, a) \neq 0 \text{ iff } s' \in \gamma(s, a)$
- cost: $S \times A \rightarrow \mathbb{R}_{\geq 0}$
 - ▶ $\text{cost}(s, a) = \text{cost of action } a \text{ in state } s$
 - ▶ may omit, default is $\text{cost}(s, a) = 1$
- Applicable(s) = $\{a \mid \gamma(s, a) \neq \emptyset\}$

Example

- $\gamma(d1, m12) = \{d2\}$
 - ▶ $\Pr(d2 | d1, m12) = 1$
- $m21, m34, m41, m43, m45, m52, m54$:
 - ▶ like $m12$
- $\gamma(d1, m14) = \{d1, d4\}$
 - ▶ $\Pr(d4 | d1, m14) = 0.5$
 - ▶ $\Pr(d1 | d1, m14) = 0.5$
- $\gamma(d2, m23) = \{d3, d5\}$
 - ▶ $\Pr(d3 | d2, m23) = 0.8$
 - ▶ $\Pr(d5 | d2, m23) = 0.2$
- there's no $m25$

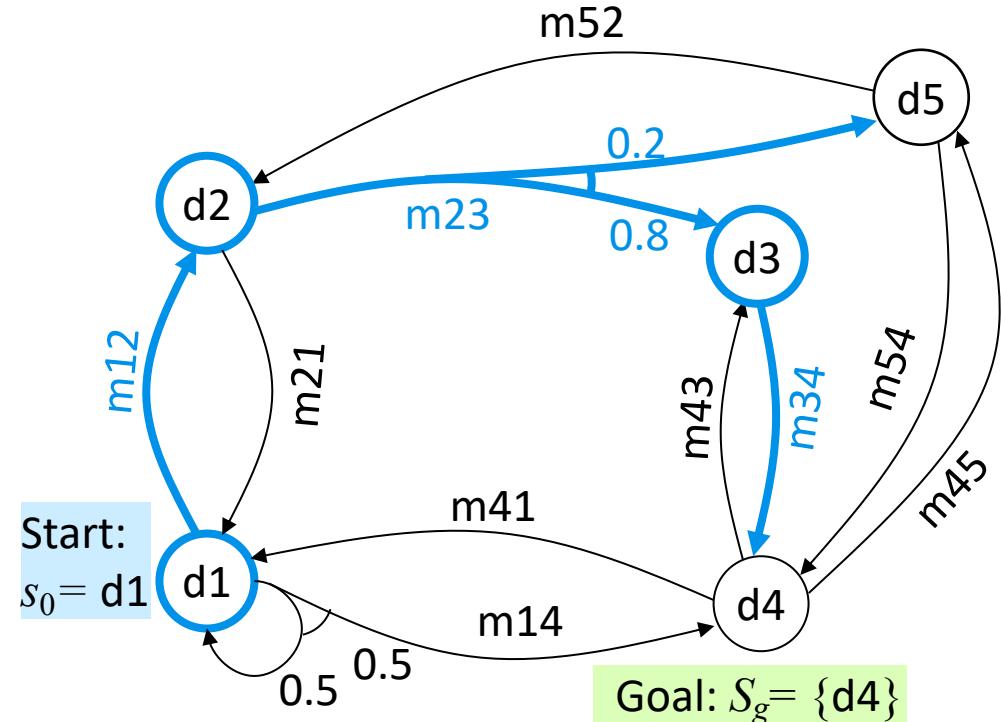


Poll: Can a plan (sequence of actions) be a solution for this problem?

1. yes
2. no

Policies, Problems

- Same as in Chapter 5:
- *Policy*
 - ▶ partial function $\pi : S \rightarrow A$ such that for every $s \in \text{Dom}(\pi) \subseteq S$, $\pi(s) \in \text{Applicable}(s)$
- *Transitive closure*
 - ▶ $\hat{\gamma}(s, \pi) = \{s\text{ and all states reachable from }s\text{ using }\pi\}$
- $\text{Graph}(s, \pi) = \text{rooted graph induced by } \pi \text{ at } s$
 - ▶ nodes: $\hat{\gamma}(s, \pi)$
 - ▶ edges: state transitions
- $\text{leaves}(s, \pi) = \hat{\gamma}(s, \pi) \setminus \text{Dom}(\pi)$



$$\pi_1 = \{(d_1, m_{12}), (d_2, m_{23}), (d_3, m_{34})\}$$

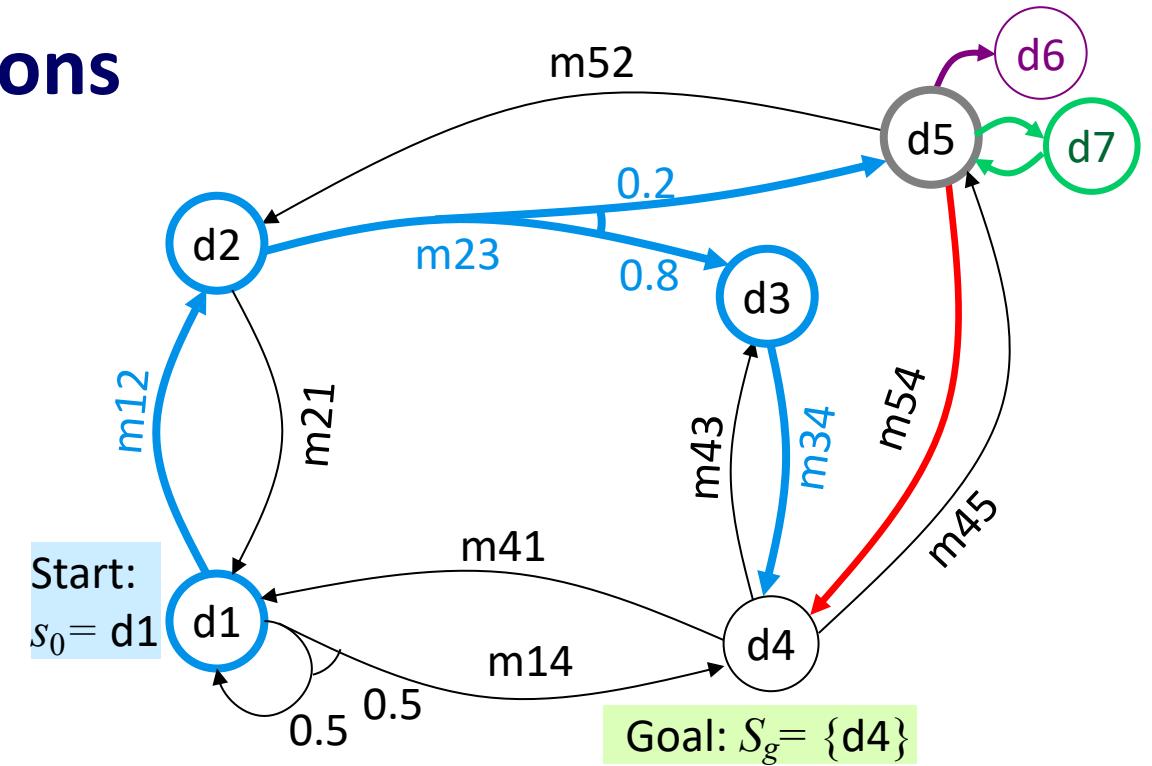
$$\text{Dom}(\pi_1) = \{d_1, d_2, d_3\}$$

$$\hat{\gamma}(d_1, \pi_1) = \{d_1, d_2, d_3, d_4, d_5\}$$

$$\begin{aligned}\text{leaves}(d_1, \pi_1) &= \hat{\gamma}(d_1, \pi_1) \setminus \text{Dom}(\pi_1) \\ &= \{d_4, d_5\}\end{aligned}$$

Solutions

- Stochastic shortest path (SSP) problem:
 - ▶ a triple (Σ, s_0, S_g)
- Solution for (Σ, s_0, S_g) :
 - ▶ A policy π for Σ such that $\hat{\gamma}((s_0, \pi) \cap S_g) \neq \emptyset$
- Unlike Chapter 5, don't require π to end at S_g
<http://www.cs.umd.edu/users/nau/apa/slides/errata.pdf>
- A solution policy π is *closed* if it doesn't stop at non-goal states unless there's no way to continue
 - ▶ for every state s in $\hat{\gamma}(s_0, \pi)$, either
 $s \in \text{Dom}(\pi)$ (i.e., $\pi(s)$ is defined)
 or $s \in S_g$
 or $\text{Applicable}(s) = \emptyset$
- For the rest of this chapter we require all solutions to be closed



$$\pi_1 = \{(d_1, m_{12}), (d_2, m_{23}), (d_3, m_{34})\}$$

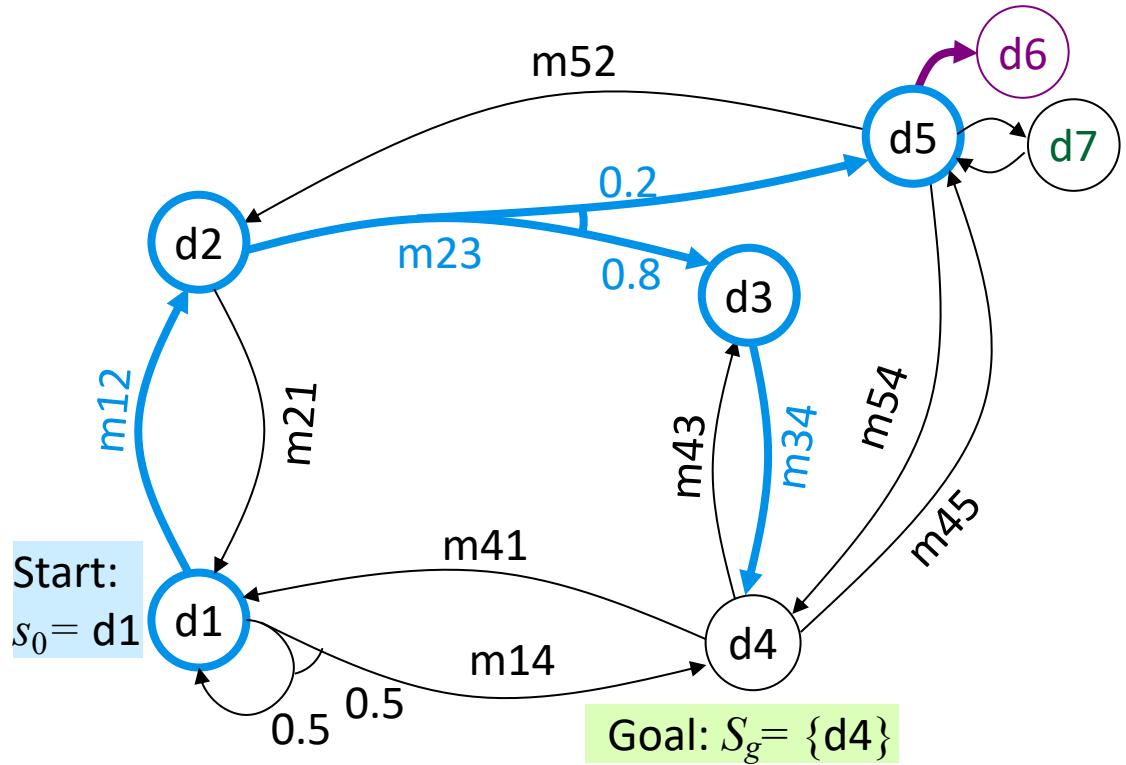
$$\pi_2 = \{(d_1, m_{12}), (d_2, m_{23}), (d_3, m_{34}), (d_5, m_{54})\}$$

$$\pi_3 = \{(d_1, m_{12}), (d_2, m_{23}), (d_3, m_{34}), (d_5, m_{56})\}$$

$$\pi_4 = \{(d_1, m_{12}), (d_2, m_{23}), (d_3, m_{34}), (d_5, m_{57}), (d_7, m_{75})\}$$

Histories

- **History:** sequence of states $\sigma = \langle s_0, s_1, s_2, \dots \rangle$
 - ▶ May be finite or infinite
 $\langle d1, d2, d3, d4 \rangle$
 $\langle d1, d2, d1, d2, \dots \rangle$
- Let $H(s, \pi) =$ set of all possible histories if we start at s and follow π
 - ▶ Stop if we reach a state s' such that $s' \notin \text{Dom}(\pi)$ or $s' \in S_g$
- If $\sigma \in H(s, \pi)$ then
 - ▶ $\Pr(\sigma | s, \pi) = \prod_{s_i, s_{i+1} \in \sigma} \Pr(s_{i+1} | s_i, \pi(s_i))$
 - ▶ Product of the probabilities of the state transitions
- $\sum_{\sigma \in H(s, \pi)} \Pr(\sigma | s, \pi) = 1$



- $\pi_3 = \{(d1, m12), (d2, m23), (d3, m34), (d5, m56)\}$
- $H(s_0, \pi_3) = \{\sigma_1, \sigma_2\}$, where:
 - ▶ $\sigma_1 = \langle d1, d2, d3, d4 \rangle$
 - ▶ $\sigma_2 = \langle d1, d2, d5, d6 \rangle$
- $\Pr(\sigma_1 | s_0, \pi_3) = 1 \times 0.8 \times 1 = 0.8$
- $\Pr(\sigma_2 | s_0, \pi_3) = 1 \times 0.2 \times 1 = 0.2$

Unsafe Solutions

- Probability of reaching a goal state:

$$\Pr(S_g | s, \pi) = \sum_{\sigma \in H(s, \pi)} \{\Pr(\sigma | s, \pi) | \sigma \text{ ends at a state in } S_g\}$$

- Equivalently:

$$\Pr(S_g | s, \pi) = \begin{cases} 1, & \text{if } s \in S_g \\ \sum_{s' \in \gamma(s, \pi(s))} \{\Pr(S_g | s', \pi)\}, & \text{otherwise} \end{cases}$$

- A solution is *unsafe* if $0 < \Pr(S_g | s_0, \pi) < 1$

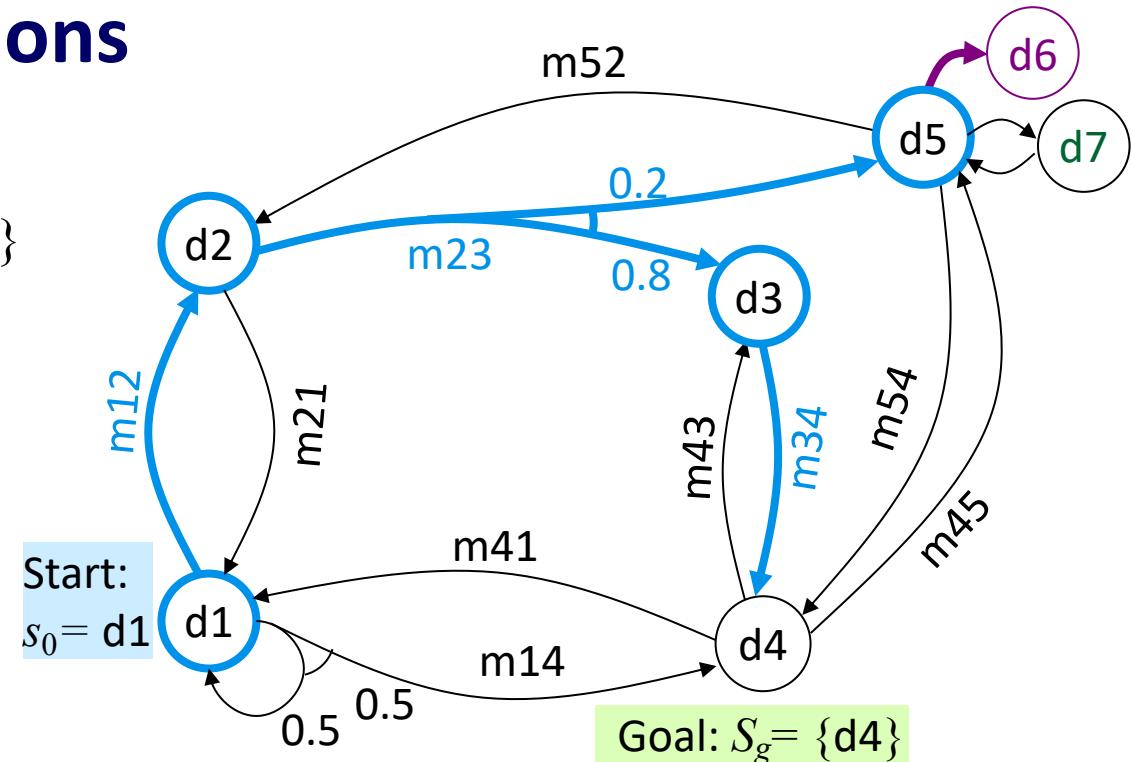
- $\pi_3 = \{(d1, m12), (d2, m23), (d3, m34), (d5, m56)\}$

- Two possible histories:

- $\sigma_1 = \langle d1, d2, d3, d4 \rangle$ ends at a goal state; $\Pr(\sigma_1 | s_0, \pi_3) = 1 \times 0.8 \times 1 = 0.8$

- $\sigma_2 = \langle d1, d2, d5, d6 \rangle$ doesn't; $\Pr(\sigma_2 | s_0, \pi_3) = 1 \times 0.2 \times 1 = 0.2$

- $\Pr(S_g | s_0, \pi_3) = \Pr(\sigma_1 | s_0, \pi_3) = 0.8$



Unsafe Solutions

- Probability of reaching a goal state:

$$\Pr(S_g | s, \pi) = \sum_{\sigma \in H(s, \pi)} \{\Pr(\sigma | s, \pi) | \sigma \text{ ends at a state in } S_g\}$$

- Equivalently:

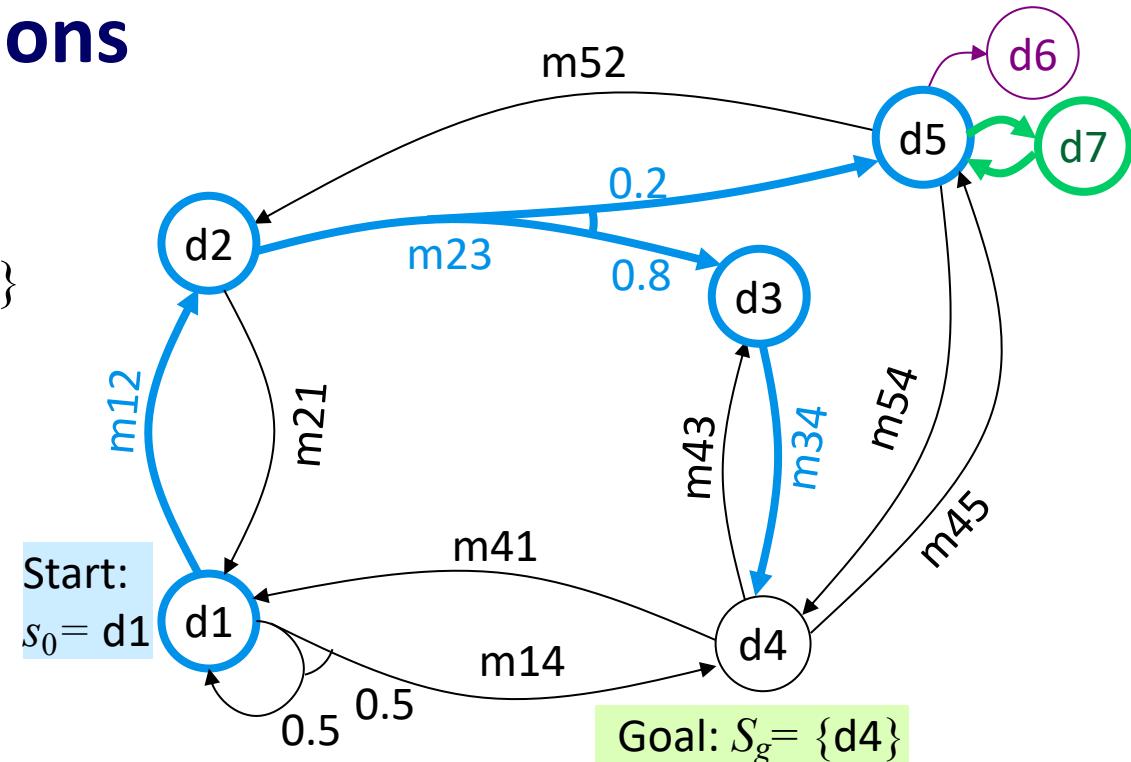
$$\Pr(S_g | s, \pi) = \begin{cases} 1, & \text{if } s \in S_g \\ \sum_{s' \in \gamma(s, \pi(s))} \{\Pr(S_g | s', \pi)\}, & \text{otherwise} \end{cases}$$

- A solution is *unsafe* if $0 < \Pr(S_g | s_0, \pi) < 1$

- $\pi_4 = \{(d1, m12), (d2, m23), (d3, m34), (d5, m57), (d7, m75)\}$

- Two possible histories

- ▶ $\sigma_1 = \langle d1, d2, d3, d4 \rangle$ ends at a goal state; $\Pr(\sigma_1 | s_0, \pi_4) = 1 \times .8 \times 1 = 0.8$
- ▶ $\sigma_3 = \langle d1, d2, d5, d7, d5, d7, \dots \rangle$ doesn't; $\Pr(\sigma_3 | s_0, \pi_4) = 1 \times .2 \times 1 \times 1 \times 1 \times \dots = 0.2$
- ▶ $\Pr(S_g | s_0, \pi_4) = \Pr(\sigma_1 | s_0, \pi_4) = 0.8$



Safe Solutions

- A solution is *safe* if $\Pr(S_g | s_0, \pi) = 1$

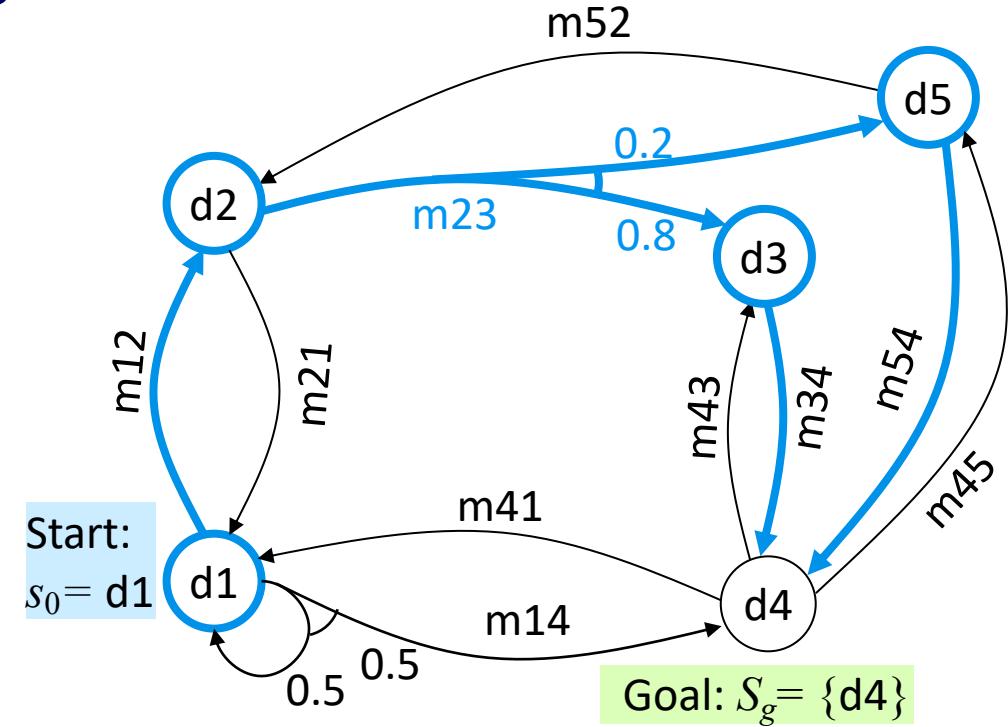
- An *acyclic* safe solution:

- ▶ $\pi_2 = \{(d1, m12), (d2, m23), (d3, m34), (d5, m54)\}$

- $H(s_0, \pi_2)$ contains two histories:

- ▶ $\sigma_1 = \langle d1, d2, d3, d4 \rangle \quad \Pr(\sigma_1 | s_0, \pi_2) = 1 \times .8 \times 1 = .8$
- ▶ $\sigma_4 = \langle d1, d2, d5, d4 \rangle \quad \Pr(\sigma_4 | s_0, \pi_2) = 1 \times .2 \times 1 = .2$

- $\Pr(S_g | s_0, \pi_2) = .8 + .2 = 1$



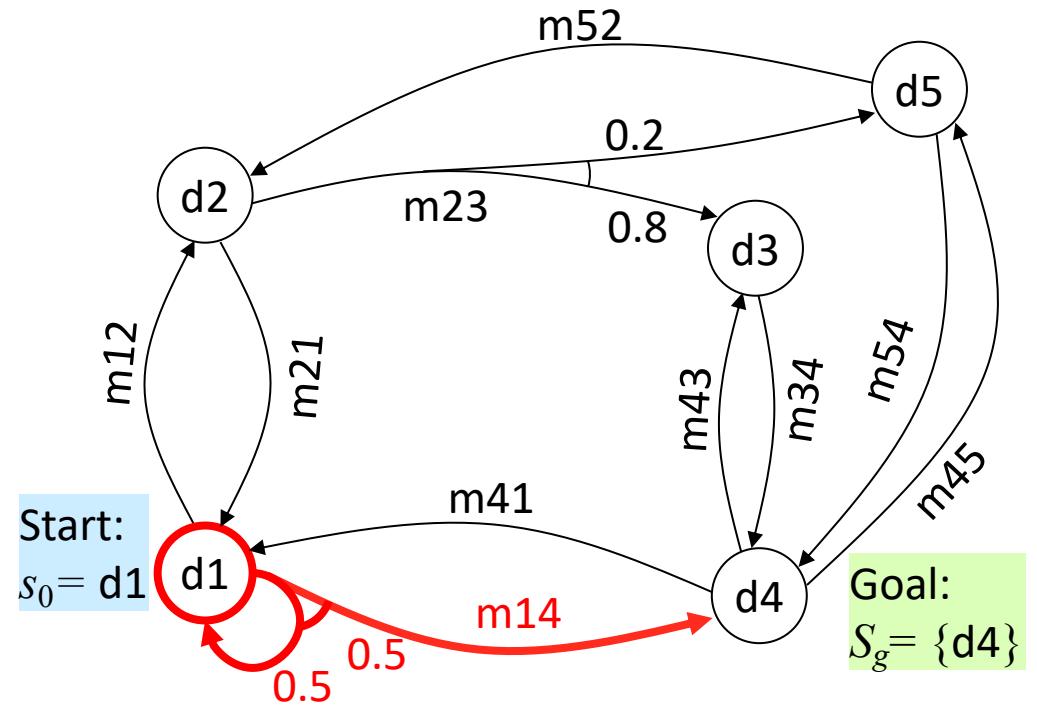
Safe Solutions

- A solution is *safe* if $\Pr(S_g | s_0, \pi) = 1$

- A *cyclic* safe solution:
 - $\pi_5 = \{(d1, m14)\}$

- $H(s_0, \pi_5)$ contains infinitely many histories:

- $\sigma_5 = \langle d1, d4 \rangle$ $\Pr(\sigma_5 | s_0, \pi_5) = \frac{1}{2}$
 - $\sigma_6 = \langle d1, d1, d4 \rangle$ $\Pr(\sigma_6 | s_0, \pi_5) = (\frac{1}{2})^2 = \frac{1}{4}$
 - $\sigma_7 = \langle d1, d1, d1, d4 \rangle$ $\Pr(\sigma_7 | s_0, \pi_5) = (\frac{1}{2})^3 = \frac{1}{8}$
 - ...
 - $\sigma_\infty = \langle d1, d1, d1, d1, d1, \dots \rangle$
- $\Pr(S_g | s_0, \pi_5) = \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \dots = 1$

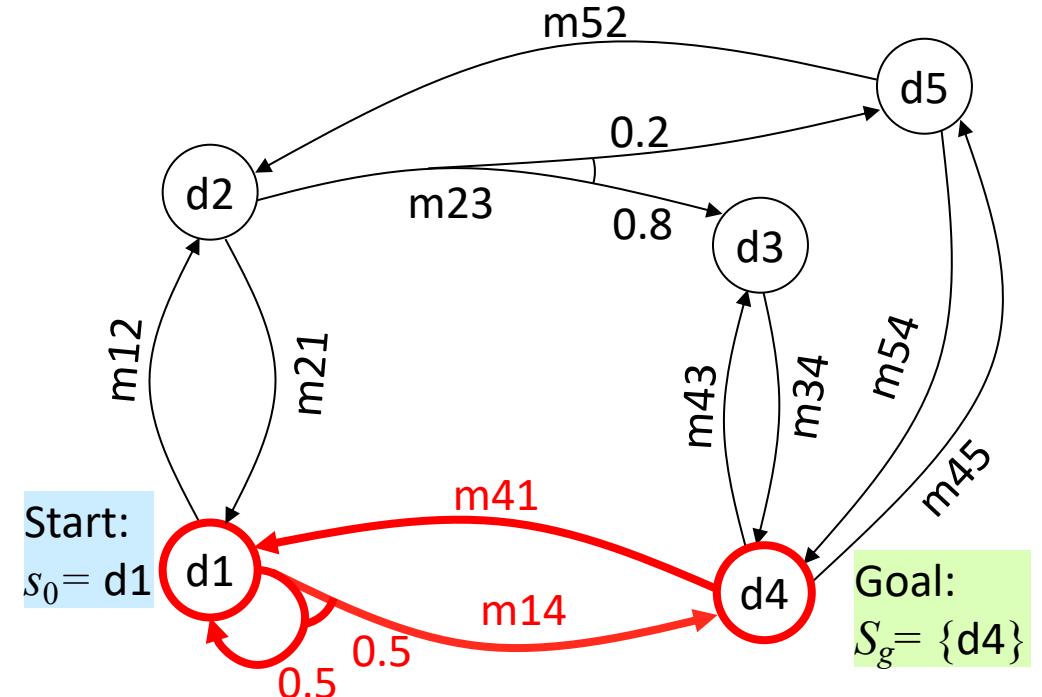


Poll: what is $\Pr(\sigma_\infty | s_0, \pi_5)$?

- 1
- 0
- a number between 0 and 1
- undefined

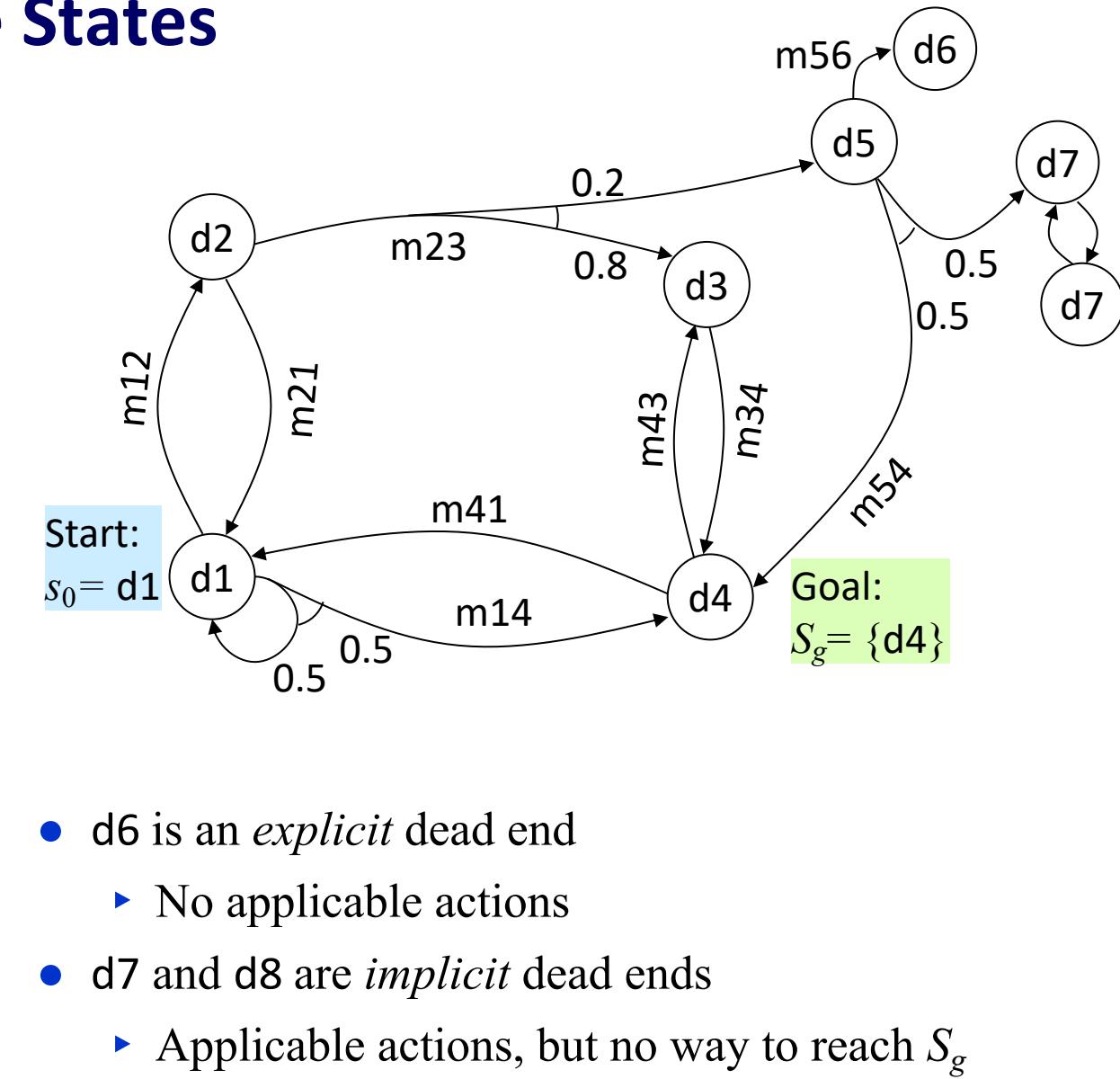
Safe Solutions

- A solution is *safe* if $\Pr(S_g | s_0, \pi) = 1$
- Another cyclic safe solution:
 $\pi_6 = \{(d1, m54), (d4, m41)\}$
- We stop when we reach a goal, so
 $(d4, m41)$ doesn't matter
 - ▶ Same histories and probabilities as for π_5



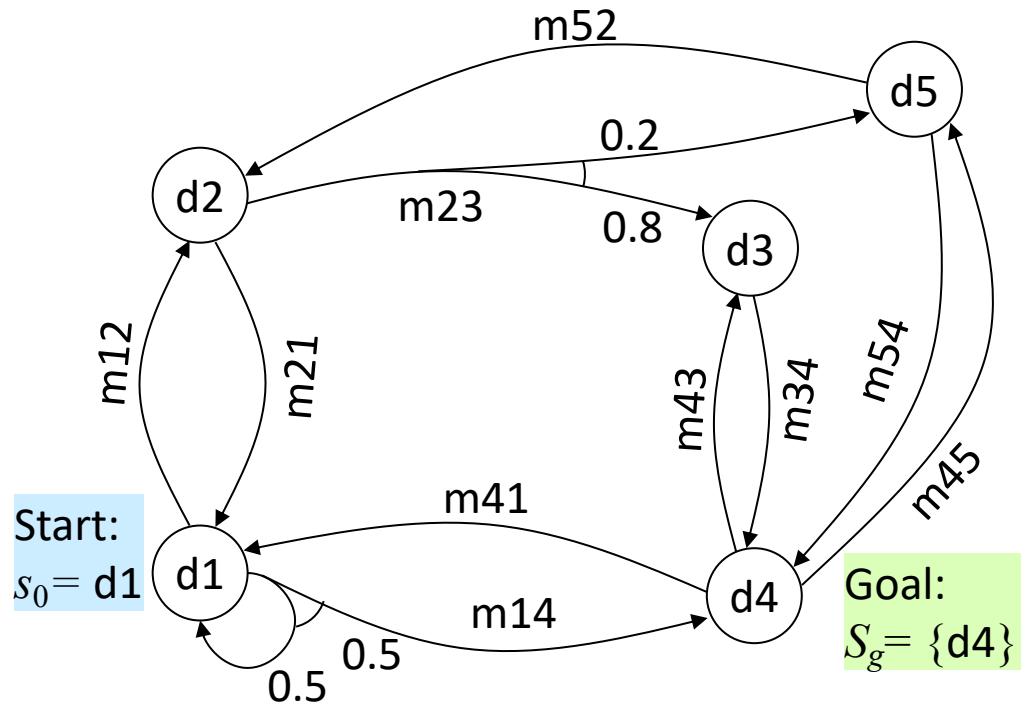
Unsafe States

- s is *safe* if $\exists \pi$ such that $\Pr(S_g | s, \pi) = 1$
 - ▶ e.g., d_3 is safe
 - ▶ same as saying (Σ, d_3, S_g) has a safe solution
- s is *unsafe* if \forall closed policy π defined at s , $0 < \Pr(S_g | s, \pi) < 1$
 - ▶ e.g., d_5 is unsafe
 - ▶ same as saying (Σ, d_5, S_g) has an unsafe solution but no safe solution
- s is a *dead end* if $\forall \pi \Pr(S_g | s, \pi) = 0$



Safe SSPs

- An SSP is *safe* if every state in the SSP is safe
 - ▶ Same as a *safely explorable* SSP in Chapter 5
- In what follows, only consider safe SSPs



Poll: Is the above SSP safe?
A. yes
B. no

- $\text{cost}(s,a) = \text{cost of using } a \text{ in } s$

- Example:

- ▶ each “horizontal” action costs 1
- ▶ each “vertical” action costs 100

- Let $\sigma = \langle s_0, s_1, s_2, \dots \rangle \in H(s_0, \pi)$

- ▶ i.e., starting at s_0 , π can produce history σ

- Then $\text{cost}(\sigma | s_0, \pi) = \sum_i \text{cost}(s_i, \pi(s_i))$

- Let π be a safe solution

- At each state $s \in \text{Dom}(\pi)$, expected cost of following π to goal:

- ▶ Weighted sum of history costs:

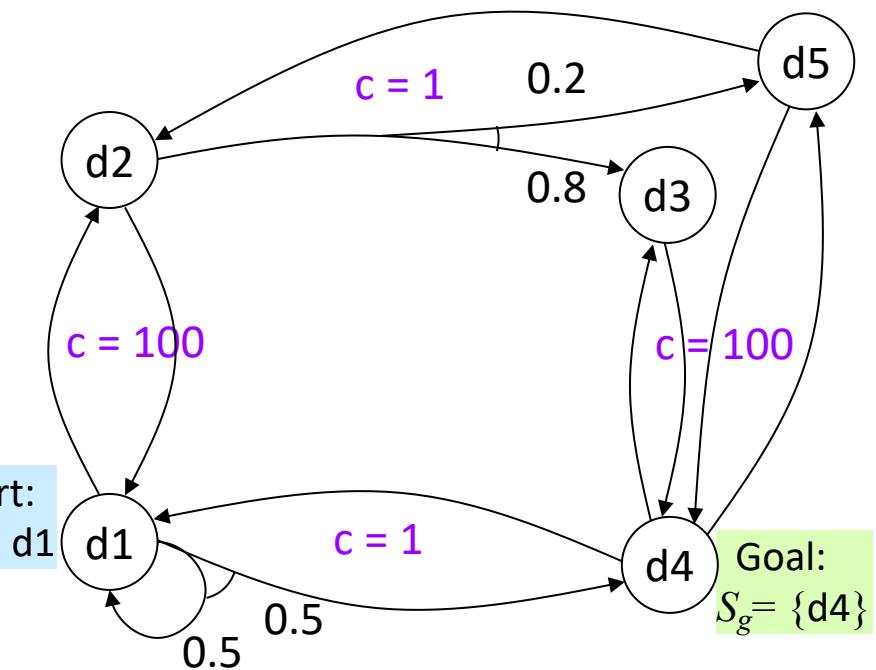
- $V^\pi(s) = \sum_{\sigma \in H(s, \pi)} \Pr(\sigma | s, \pi) \text{cost}(\sigma | s, \pi)$

My version

- ▶ Recursive equation

$$V^\pi(s) = \begin{cases} 0, & \text{if } s \in S_g \\ \text{cost}(s, \pi(s)) + \sum_{s' \in \gamma(s, \pi(s))} \Pr(s' | s, \pi(s)) V^\pi(s'), & \text{otherwise} \end{cases}$$

From the book

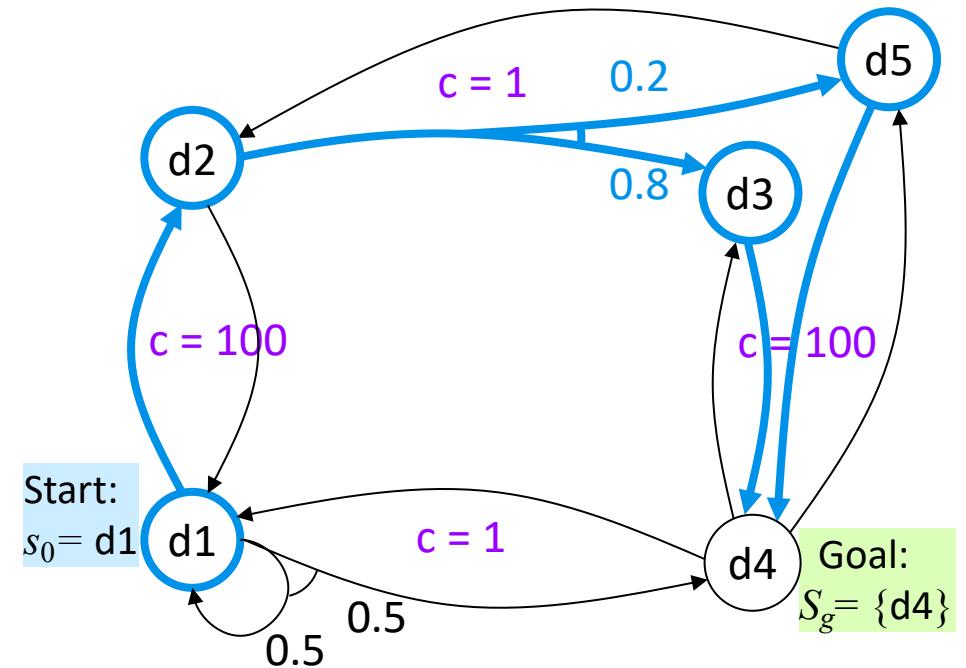


Poll: Are the two equations equivalent?

- A. yes
B. no

Example

- $\pi_3 = \{(d1, m12), (d2, m23), (d3, m34), (d5, m54)\}$
- Weighted sum of history costs:
 - $\sigma_1 = \langle d1, d2, d3, d4 \rangle$
 - $\Pr(\sigma_1 | s_0, \pi_3) = 0.8$
 - $\text{cost}(\sigma_1 | s_0, \pi_3) = 100 + 1 + 100 = 201$
 - $\sigma_2 = \langle d1, d2, d5, d4 \rangle$
 - $\Pr(\sigma_2 | s_0, \pi_3) = 0.2$
 - $\text{cost}(\sigma_2 | s_0, \pi_3) = 100 + 1 + 100 = 201$
- $V^{\pi_3}(d1) = .8(201) + .2(201) = 201$



- Recursive equation \Rightarrow 4 equations, 4 unknowns

$$V^{\pi_3}(d1) = 100 + V^{\pi_3}(d2)$$

$$V^{\pi_3}(d2) = 1 + .8(V^{\pi_3}(d3)) + .2(V^{\pi_3}(d5))$$

$$V^{\pi_3}(d3) = 100 + V^{\pi_3}(d4)$$

$$V^{\pi_3}(d5) = 100 + V^{\pi_3}(d4)$$

$$V^{\pi_3}(d4) = 0$$
- So $V^{\pi_3}(d1) = 100 + 1 + .8(100) + .2(100) = 201$

Example

- $\pi_7 = \{(d1, m14), (d2, m23), (d3, m34), (d5, m54)\}$

- Weighted sum of history costs:

- ▶ $\sigma_5 = \langle d1, d4 \rangle$

$$\Pr(\sigma_5 | \pi_7) = \frac{1}{2}, \quad \text{cost}(\sigma_5 | \pi_7) = 1$$

- ▶ $\sigma_6 = \langle d1, d1, d4 \rangle$

$$\Pr(\sigma_6 | \pi_7) = (\frac{1}{2})^2, \quad \text{cost}(\sigma_6 | \pi_7) = 2$$

- ▶ $\sigma_7 = \langle d1, d1, d1, d4 \rangle$

$$\Pr(\sigma_7 | \pi_7) = (\frac{1}{2})^3, \quad \text{cost}(\sigma_7 | \pi_7) = 3$$

...

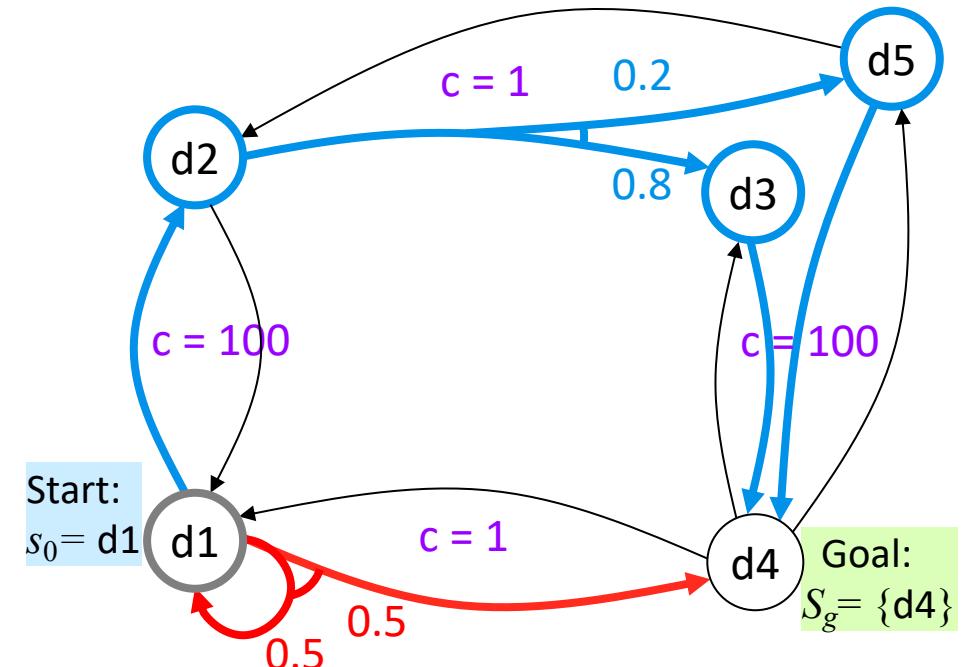
- $V^{\pi_7}(d1) = (\frac{1}{2})1 + (\frac{1}{2})^2 2 + (\frac{1}{2})^3 3 + \dots = 2$

- Recursive equation:

$$V^{\pi_7}(d1) = 1 + \frac{1}{2}(0) + \frac{1}{2}(V^{\pi_7}(d1))$$

$$\frac{1}{2}V^{\pi_7}(d1) = 1$$

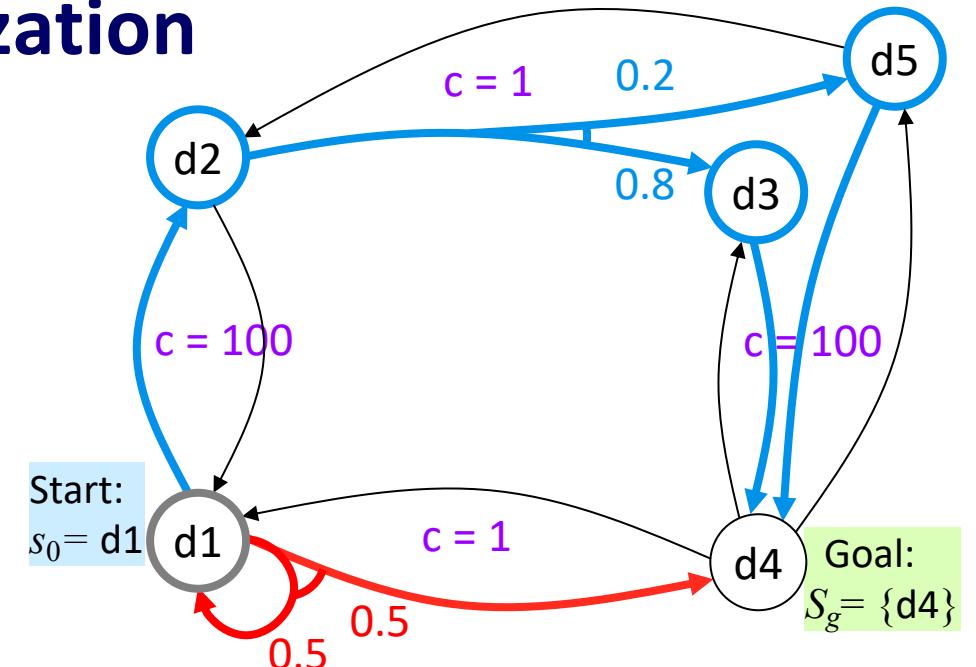
$$V^{\pi_7}(d1) = 2$$



- Given safe solution π ,
 - ▶ Compute V^π by solving n linear equations, n unknowns
 - ▶ $n = \text{number of states reachable from } s_0 \text{ using } \pi = |\hat{\gamma}(s_0, \pi)|$

Planning as Optimization

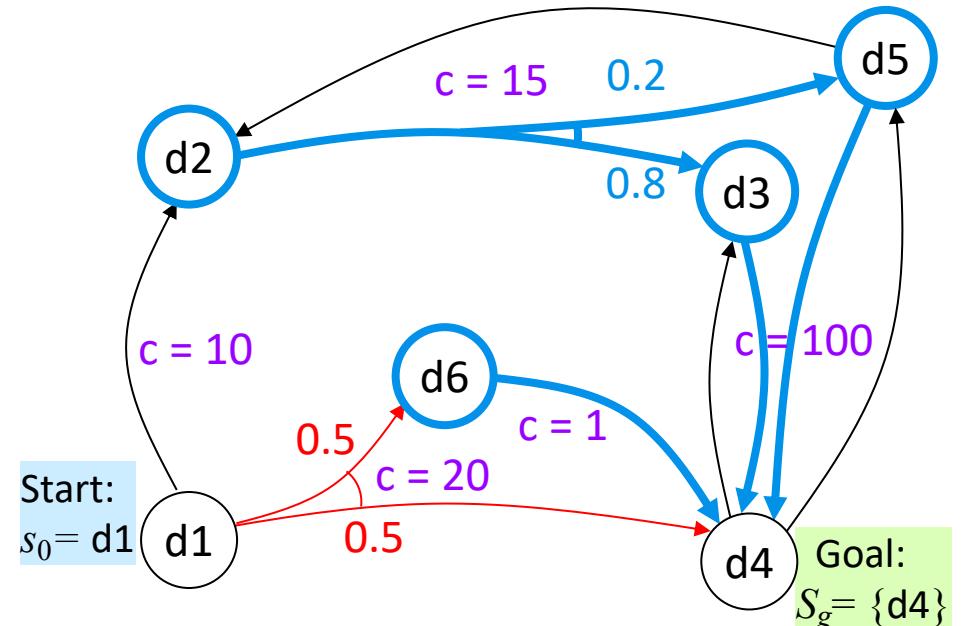
- Let π and π' be safe solutions
 - π dominates π' if $V^\pi(s) \leq V^{\pi'}(s)$ at every state s where they're both defined
 - i.e., every state $s \in \text{Dom}(\pi) \cap \text{Dom}(\pi')$
- On the previous two slides
 - $\pi_3 = \{(d1, m12), (d2, m23), (d3, m34), (d5, m54)\}$
 - $\pi_7 = \{(d1, m14), (d2, m23), (d3, m34), (d5, m54)\}$
 - $d1$ is the only state where they differ
 - $V^{\pi_3}(d1) = 201; V^{\pi_7}(d1) = 2$
 - π_7 dominates π_3
- On slide 11, $\pi_5 = \{(d1, m14)\}$
 - $d1$ is the only state where they're both defined
 - $V^{\pi_3}(d1) = 201; V^{\pi_5}(d1) = 2$
 - π_5 dominates π_3



- π is optimal if π dominates every safe solution
- If π and π' are both optimal, then $V^\pi(s) = V^{\pi'}(s)$ at every state where they're both defined
 - e.g., π_5 and π_7

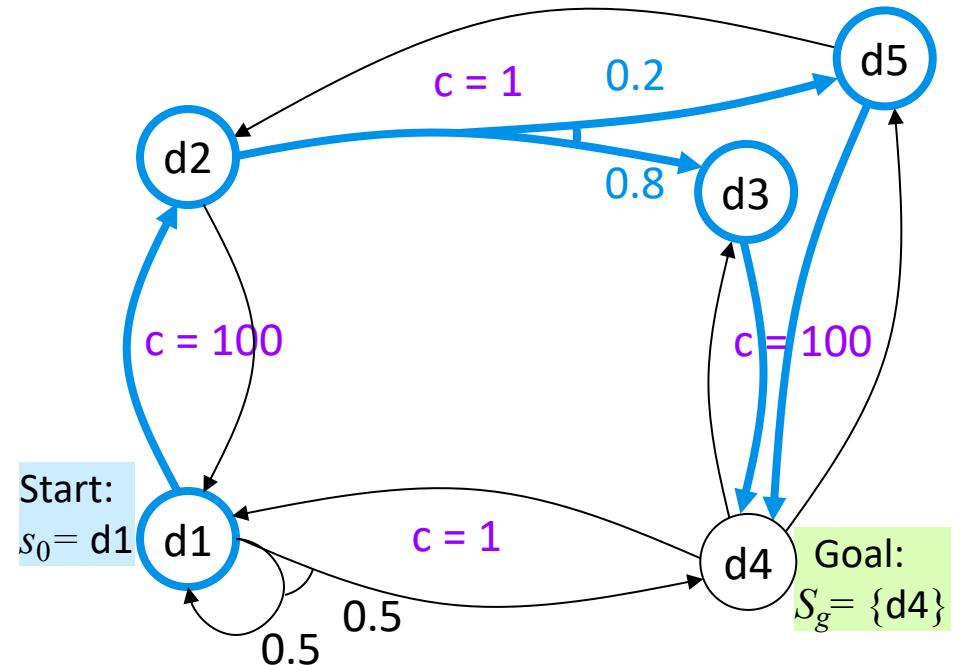
Planning as Optimization

- Let $V^*(s)$ = expected cost of an optimal safe solution
- Example:
 - At state d_1 , applicable actions m_{12} and m_{14}
 - Suppose we know $V^*(d_2)$, $V^*(d_4)$, $V^*(d_6)$
- Then
 - $\text{cost}(d_1, m_{12}) = 10 + V^*(d_2)$
 - $\text{cost}(d_1, m_{14}) = 20 + 0.5 V^*(d_6) + 0.5 V^*(d_4)$
 - $V^*(d_1) = \min(\text{cost}(d_1, m_{12}), \text{cost}(d_1, m_{14}))$
- In general,
$$V^*(s) = \begin{cases} 0, & \text{if } s \text{ is a goal} \\ \min_{a \in \text{Applicable}(s)} \{ \text{cost}(s, a) + \sum_{s' \in \gamma(s, a)} \Pr(s' | s, a) V^*(s') \}, & \text{otherwise} \end{cases}$$
- Optimality principle* (Bellman's theorem)



Cost to Go

- Let (Σ, s_0, S_g) be a *safe* SSP
 - i.e., S_g is reachable from every state
 - same as *safely explorable* in Chapter 5
- Let π be a safe solution that's defined at all non-goal states
 - i.e., $\text{Dom}(\pi) = S \setminus S_g$
- Compute V^π as $|S|$ equations, $|S|$ unknowns
- Let $s \in S$, $a \in \text{Applicable}(s)$
 - Cost-to-go*:
 - expected cost at s if we first use a , then use π afterward
 - $Q^\pi(s, a) = \text{cost}(s, a) + \sum_{s' \in \gamma(s, a')} \Pr(s' | s, a') V^\pi(s')$

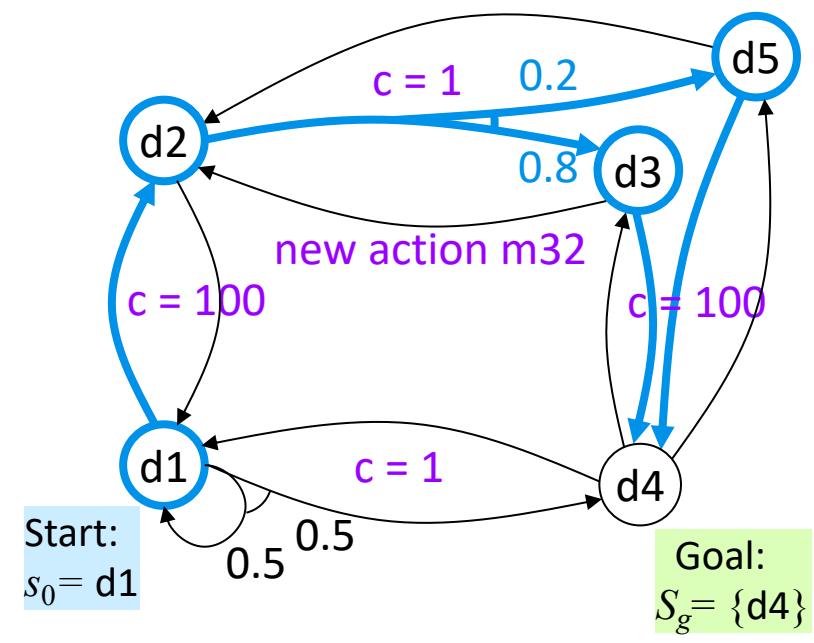


- For every $s \in S \setminus S_g$
let $\pi'(s) \in \operatorname{argmin}_{a' \in \text{Applicable}(s)} Q^\pi(s, a')$

Poll: Does π' dominate π ?

- always
- sometimes
- never

Example



$$\pi = \{(d_1, m_{12}), (d_2, m_{23}), (d_3, m_{34}), (d_5, m_{54})\}$$

$$V^\pi(d_4) = 0$$

$$V^\pi(d_3) = 100 + V^\pi(d_4) = 100$$

$$V^\pi(d_5) = 100 + V^\pi(d_4) = 100$$

$$V^\pi(d_2) = 1 + (0.8 V^\pi(d_3) + 0.2 V^\pi(d_5)) = 101$$

$$V^\pi(d_1) = 100 + V^\pi(d_2) = 201$$

$$Q^\pi(d_1, m_{12}) = 100 + 101 = 201$$

$$Q^\pi(d_1, m_{14}) = 1 + \frac{1}{2}(201) + \frac{1}{2}(0) = 101.5 \\ \Rightarrow \operatorname{argmin}_a Q^\pi(d_1, a) = m_{14}$$

$$Q^\pi(d_2, m_{23}) = 1 + (0.8(100) + 0.2(100)) = 101$$

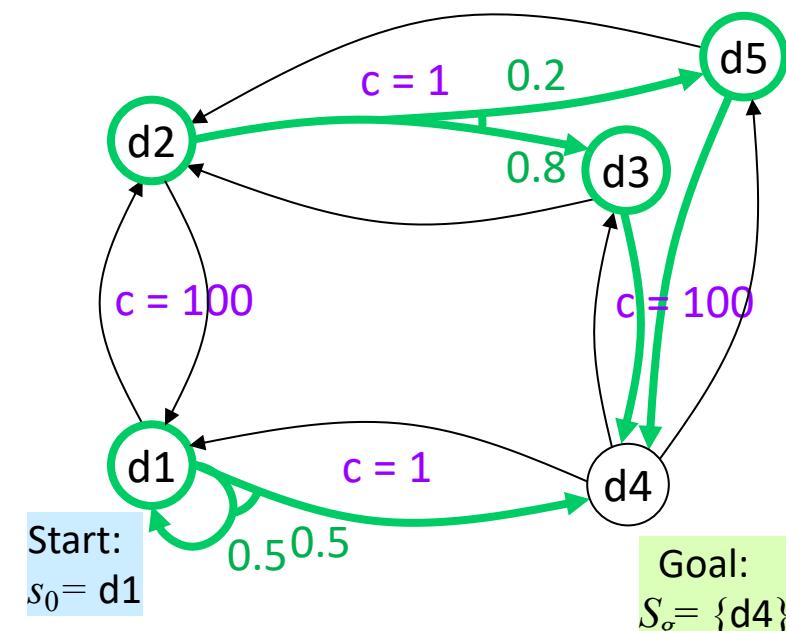
$$Q^\pi(d_2, m_{21}) = 100 + 201 = 301 \\ \Rightarrow \operatorname{argmin}_a Q^\pi(d_2, a) = m_{23}$$

$$Q^\pi(d_3, m_{34}) = 100 + 0 = 100$$

$$Q^\pi(d_3, m_{32}) = 1 + 101 = 102 \\ \Rightarrow \operatorname{argmin}_a Q^\pi(d_3, a) = m_{34}$$

$$Q^\pi(d_5, m_{54}) = 100 + 0 = 100$$

$$Q^\pi(d_5, m_{52}) = 1 + 101 = 102 \\ \Rightarrow \operatorname{argmin}_a Q^\pi(d_5, a) = m_{54}$$



$$\pi' = \{(d_1, m_{14}), (d_2, m_{23}), (d_3, m_{34}), (d_5, m_{54})\}$$

$$V^{\pi'}(d_4) = 0$$

$$V^{\pi'}(d_3) = 100 + V^{\pi'}(d_4) = 100$$

$$V^{\pi'}(d_5) = 100 + V^{\pi'}(d_4) = 100$$

$$V^{\pi'}(d_2) = 1 + (0.8 V^{\pi'}(d_3) + 0.2 V^{\pi'}(d_5)) = 101$$

$$V^{\pi'}(d_1) = 1 + \frac{1}{2}V^{\pi'}(d_1) + \frac{1}{2}V^{\pi'}(d_4) \\ \Rightarrow V^{\pi'}(d_1) = 2$$

- $\text{PI}(\Sigma, s_0, S_g, \pi_0)$

$$\pi \leftarrow \pi_0$$

loop

compute $\{V^\pi(s) \mid s \in S\} \leftarrow |S| \text{ equations, } |S| \text{ unknowns}$

for every $s \in S \setminus S_g$ do

$\pi'(s) \leftarrow \operatorname{argmin}_{a \in \text{Applicable}(s)} Q^\pi(s, a) \leftarrow E(\text{cost of using } a \text{ then } \pi)$

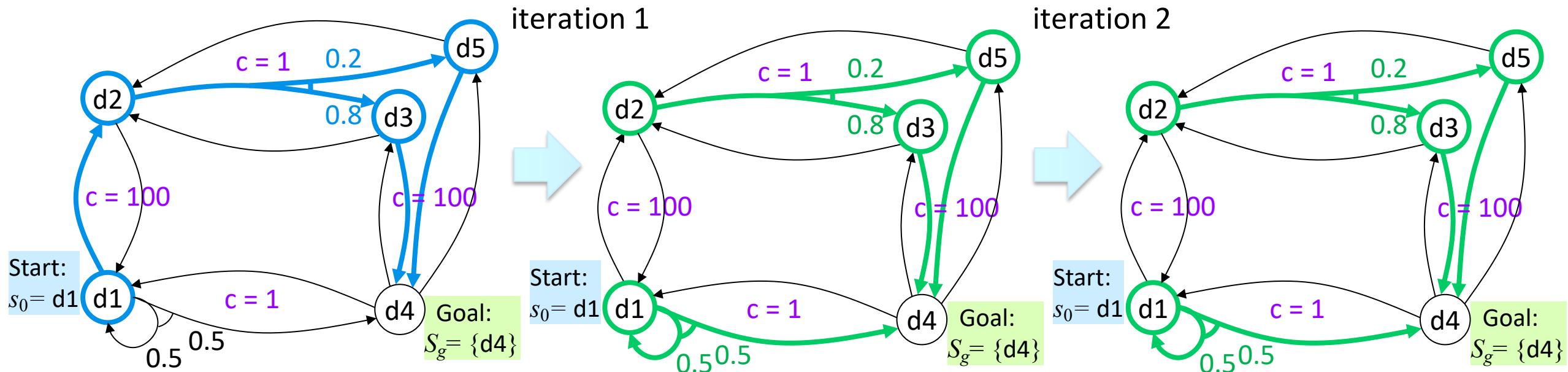
if $\pi' = \pi$ then

return π

$\pi \leftarrow \pi'$

Policy Iteration

- Converges in a finite number of iterations
- Example:



Value Iteration

- V_0 : a heuristic function
 - e.g., adapt an h heuristic from Chapter 2
 - ▶ $V_0(s) \approx$ optimal expected cost of getting from s to S_g
 - ▶ Require $V_0(s) = 0$ for every $s \in S_g$
- $\eta > 0$: for testing approximate convergence
- V (global variable): an array
 - ▶ $V(s) \approx$ optimal expected cost of getting from s to S_g
 - ▶ Initially $V(s) = V_0(s) \forall s$
 - ▶ Improve the estimates as we go along
- π (global variable): policy computed from V
- Difference from the book:
 - ▶ In the book, VI computes r as a separate step, not in Bellman-Update

```
VI( $\Sigma, s_0, S_g, V_0, \eta$ )
  global  $V \leftarrow V_0$ 
  global  $\pi \leftarrow \emptyset$ 
  loop
     $r \leftarrow \max \{ \text{Bellman-Update}(s) \mid s \in S \setminus S_g \}$ 
    if  $r \leq \eta$  then return  $\pi$ 
```

```
Bellman-Update( $s$ )
  global  $V, \pi$ 
   $v_{\text{old}} \leftarrow V(s)$ 
  for every  $a \in \text{Applicable}(s)$  do
     $Q(s,a) \leftarrow \text{cost}(s,a) + \sum_{s' \in S} \Pr(s'|s,a) V(s')$ 
     $V(s) \leftarrow \min_{a \in \text{Applicable}(s)} Q(s,a)$ 
     $\pi(s) \leftarrow \operatorname{argmin}_{a \in \text{Applicable}(s)} Q(s,a)$ 
  return  $|V(s) - v_{\text{old}}|$ 
```

Iteration 1

```

VI( $\Sigma, s_0, S_g, V_0, \eta$ )
global  $V \leftarrow V_0$ 
global  $\pi \leftarrow \emptyset$ 
loop
   $r \leftarrow \max \{ \text{Bellman-Update}(s) \mid s \in S \setminus S_g \}$ 
  if  $r \leq \eta$  then return  $\pi$ 

```

Bellman-Update(s)

global V, π

$v_{\text{old}} \leftarrow V(s)$

for every $a \in \text{Applicable}(s)$ do

$$Q(s, a) \leftarrow \text{cost}(s, a) + \sum_{s' \in S} \Pr(s'|s, a) V(s')$$

$$V(s) \leftarrow \min_{a \in \text{Applicable}(s)} Q(s, a)$$

$$\pi(s) \leftarrow \operatorname{argmin}_{a \in \text{Applicable}(s)} Q(s, a)$$

return $|V(s) - v_{\text{old}}|$

$$\eta = 0.2$$

initial values

$V(d1) = 0$
$V(d2) = 0$
$V(d3) = 0$
$V(d5) = 0$

$$Q(d1, m12) = 100 + 0 = 100$$

$$Q(d1, m14) = 1 + \frac{1}{2}(0) + \frac{1}{2}(0) = 1$$

$$V(d1) = 1; \pi(d1) = m14; \quad 1 - 0 = 1$$

$$Q(d2, m21) = 100 + 1 = 101$$

$$Q(d2, m23) = 1 + .8(0) + .2(0) = 1$$

$$V(d2) = 1; \pi(d2) = m23; \quad 1 - 0 = 1$$

$$Q(d3, m32) = 1 + 1 = 2$$

$$Q(d3, m34) = 100 + 0 = 100$$

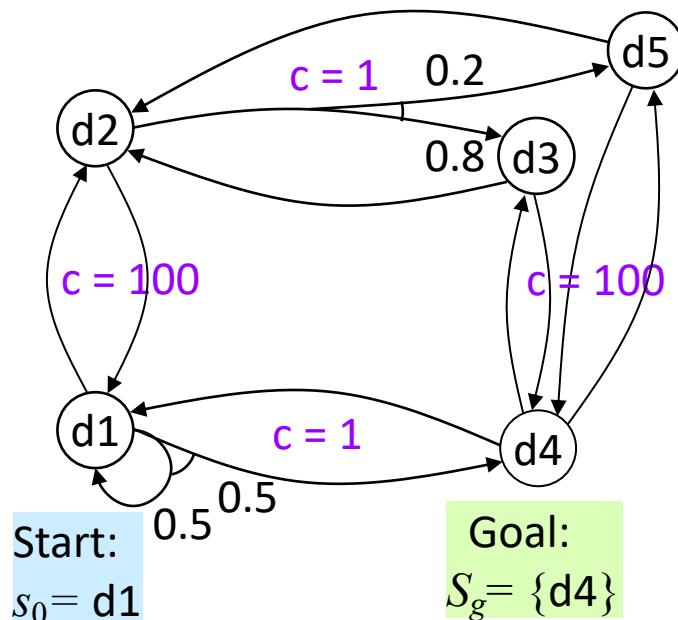
$$V(d3) = 2; \pi(d3) = m32; \quad 2 - 0 = 2$$

$$Q(d5, m52) = 1 + 1 = 2$$

$$Q(d5, m54) = 100 + 0 = 100$$

$$V(d5) = 2; \pi(d5) = m52; \quad 2 - 0 = 2$$

$$r = \max(1, 1, 2, 2) = 2$$



Iteration 2

```

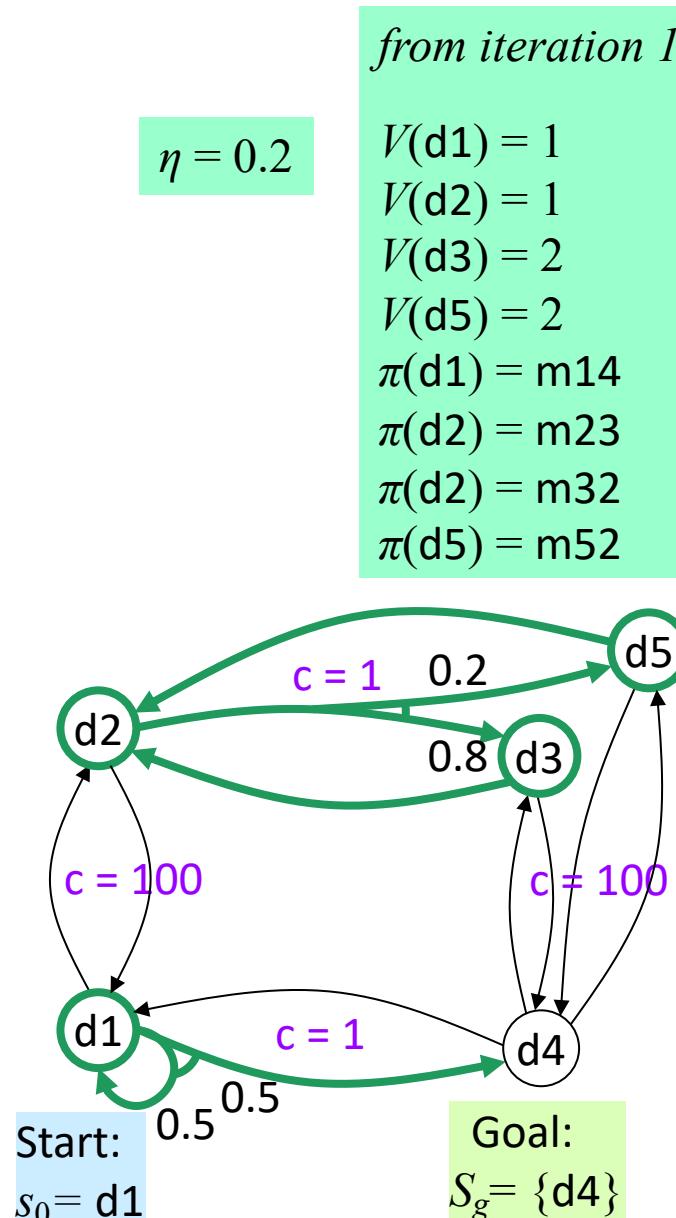
VI( $\Sigma, s_0, S_g, V_0, \eta$ )
global  $V \leftarrow V_0$ 
global  $\pi \leftarrow \emptyset$ 
loop
   $r \leftarrow \max \{ \text{Bellman-Update}(s) \mid s \in S \setminus S_g \}$ 
  if  $r \leq \eta$  then return  $\pi$ 

```

```

Bellman-Update( $s$ )
global  $V, \pi$ 
 $v_{\text{old}} \leftarrow V(s)$ 
for every  $a \in \text{Applicable}(s)$  do
   $Q(s,a) \leftarrow \text{cost}(s,a) + \sum_{s' \in S} \Pr(s'|s,a) V(s')$ 
 $V(s) \leftarrow \min_{a \in \text{Applicable}(s)} Q(s,a)$ 
 $\pi(s) \leftarrow \operatorname{argmin}_{a \in \text{Applicable}(s)} Q(s,a)$ 
return  $|V(s) - v_{\text{old}}|$ 

```



$$Q(d1,m12) = 100 + 1 = 101$$

$$Q(d1,m14) = 1 + \frac{1}{2}(1) + \frac{1}{2}(0) = 1\frac{1}{2}$$

$$V(d1) = 1\frac{1}{2}; \pi(d1) = m14; 1\frac{1}{2} - 1 = \frac{1}{2}$$

$$Q(d2,m21) = 100 + 1\frac{1}{2} = 101\frac{1}{2}$$

$$Q(d2,m23) = 1 + .8(2) + .2(2) = 3$$

$$V(d2) = 3; \pi(d2) = m23; 3 - 1 = 2$$

$$Q(d3,m32) = 1 + 3 = 4$$

$$Q(d3,m34) = 100 + 0 = 100$$

$$V(d3) = 4; \pi(d3) = m32; 4 - 2 = 2$$

$$Q(d5,m52) = 1 + 3 = 4$$

$$Q(d5,m54) = 100 + 0 = 100$$

$$V(d5) = 4; \pi(d5) = m52; 4 - 2 = 2$$

$$r = \max(\frac{1}{2}, 2, 2, 2) = 2$$

Iteration 3

```

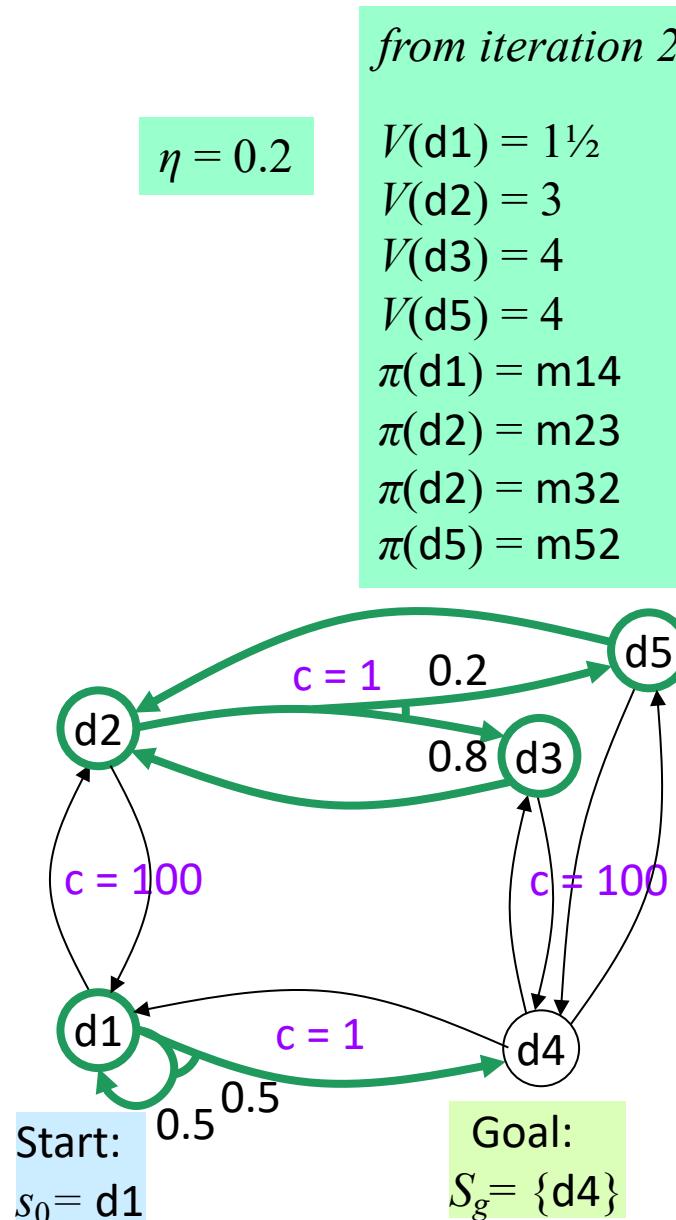
VI( $\Sigma, s_0, S_g, V_0, \eta$ )
global  $V \leftarrow V_0$ 
global  $\pi \leftarrow \emptyset$ 
loop
   $r \leftarrow \max \{ \text{Bellman-Update}(s) \mid s \in S \setminus S_g \}$ 
  if  $r \leq \eta$  then return  $\pi$ 

```

```

Bellman-Update( $s$ )
global  $V, \pi$ 
 $v_{\text{old}} \leftarrow V(s)$ 
for every  $a \in \text{Applicable}(s)$  do
   $Q(s,a) \leftarrow \text{cost}(s,a) + \sum_{s' \in S} \Pr(s'|s,a) V(s')$ 
 $V(s) \leftarrow \min_{a \in \text{Applicable}(s)} Q(s,a)$ 
 $\pi(s) \leftarrow \operatorname{argmin}_{a \in \text{Applicable}(s)} Q(s,a)$ 
return  $|V(s) - v_{\text{old}}|$ 

```



$$Q(d1,m12) = 100 + 3 = 103$$

$$Q(d1,m14) = 1 + \frac{1}{2}(1\frac{1}{2}) + \frac{1}{2}(0) = 1\frac{3}{4}$$

$$V(d1) = 1\frac{3}{4}; \pi(d1) = m14; 1\frac{3}{4} - 1\frac{1}{2} = \frac{1}{4}$$

$$Q(d2,m21) = 100 + 1\frac{3}{4} = 101\frac{3}{4}$$

$$Q(d2,m23) = 1 + .8(4) + .2(4) = 5$$

$$V(d2) = 5; \pi(d2) = m23; 5 - 3 = 2$$

$$Q(d3,m32) = 1 + 5 = 6$$

$$Q(d3,m34) = 100 + 0 = 100$$

$$V(d3) = 6; \pi(d3) = m32; 6 - 4 = 2$$

$$Q(d5,m52) = 1 + 5 = 6$$

$$Q(d5,m54) = 100 + 0 = 100$$

$$V(d5) = 6; \pi(d5) = m52; 6 - 4 = 2$$

$$r = \max(\frac{1}{4}, 2, 2, 2) = 2$$

Iteration 4

```

VI( $\Sigma, s_0, S_g, V_0, \eta$ )
global  $V \leftarrow V_0$ 
global  $\pi \leftarrow \emptyset$ 
loop
   $r \leftarrow \max \{ \text{Bellman-Update}(s) \mid s \in S \setminus S_g \}$ 
  if  $r \leq \eta$  then return  $\pi$ 

```

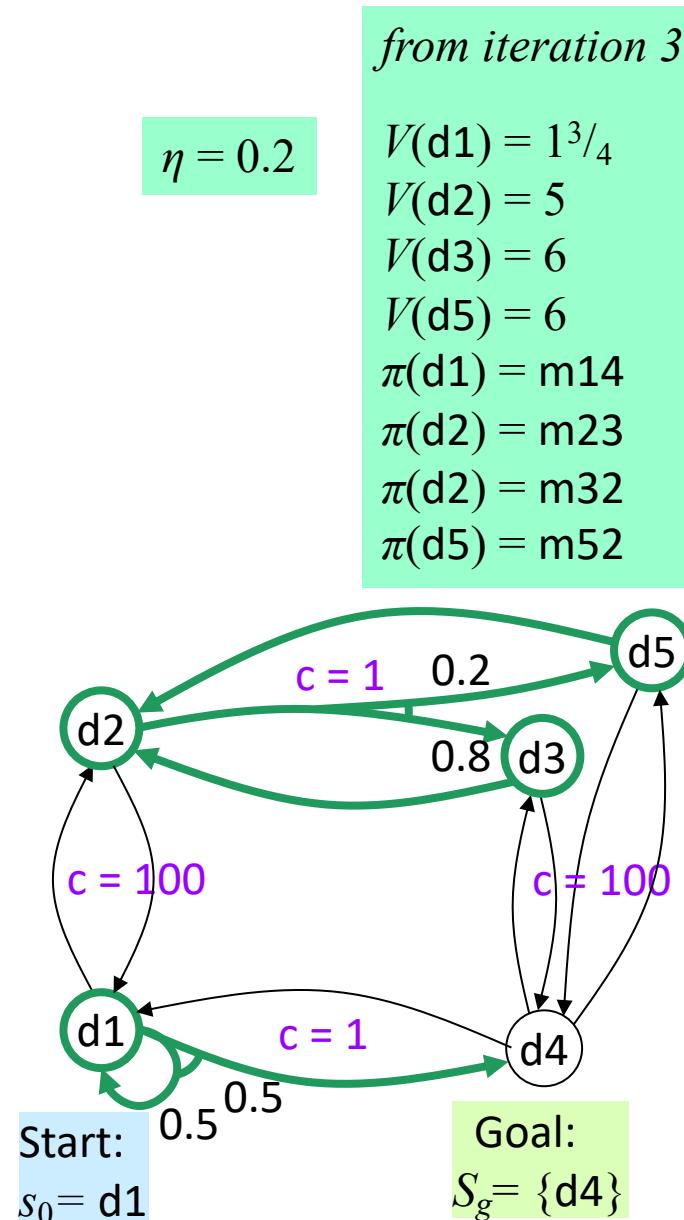
```

Bellman-Update( $s$ )
global  $V, \pi$ 
 $v_{\text{old}} \leftarrow V(s)$ 
for every  $a \in \text{Applicable}(s)$  do
   $Q(s,a) \leftarrow \text{cost}(s,a) + \sum_{s' \in S} \Pr(s'|s,a) V(s')$ 
   $V(s) \leftarrow \min_{a \in \text{Applicable}(s)} Q(s,a)$ 
   $\pi(s) \leftarrow \operatorname{argmin}_{a \in \text{Applicable}(s)} Q(s,a)$ 
return  $|V(s) - v_{\text{old}}|$ 

```

Poll: How many iterations until termination?

A. $i \leq 10$	D. $40 < i \leq 80$
B. $10 < i \leq 20$	E. $80 < i \leq 160$
C. $20 < i \leq 40$	F. $160 < i$



$$Q(d1,m12) = 100 + 5 = 105$$

$$Q(d1,m14) = 1 + \frac{1}{2}(1^{3/4}) + \frac{1}{2}(0) = 1^{7/8}$$

$$V(d1) = 1^{7/8}; \pi(d1) = m14; 1^{7/8} - 1^{3/4} = 1/8$$

$$Q(d2,m21) = 100 + 1^{7/8} = 101^{7/8}$$

$$Q(d2,m23) = 1 + .8(6) + .2(6) = 7$$

$$V(d2) = 7; \pi(d2) = m23; 7 - 5 = 2$$

$$Q(d3,m32) = 1 + 7 = 8$$

$$Q(d3,m34) = 100 + 0 = 100$$

$$V(d3) = 8; \pi(d3) = m32; 8 - 6 = 2$$

$$Q(d5,m52) = 1 + 7 = 8$$

$$Q(d5,m54) = 100 + 0 = 100$$

$$V(d5) = 8; \pi(d5) = m52; 8 - 6 = 2$$

$$r = \max(1/8, 2, 2, 2) = 2$$

Iteration 1, with a better V_0

```

VI( $\Sigma, s_0, S_g, V_0, \eta$ )
global  $V \leftarrow V_0$ 
global  $\pi \leftarrow \emptyset$ 
loop
   $r \leftarrow \max \{ \text{Bellman-Update}(s) \mid s \in S \setminus S_g \}$ 
  if  $r \leq \eta$  then return  $\pi$ 

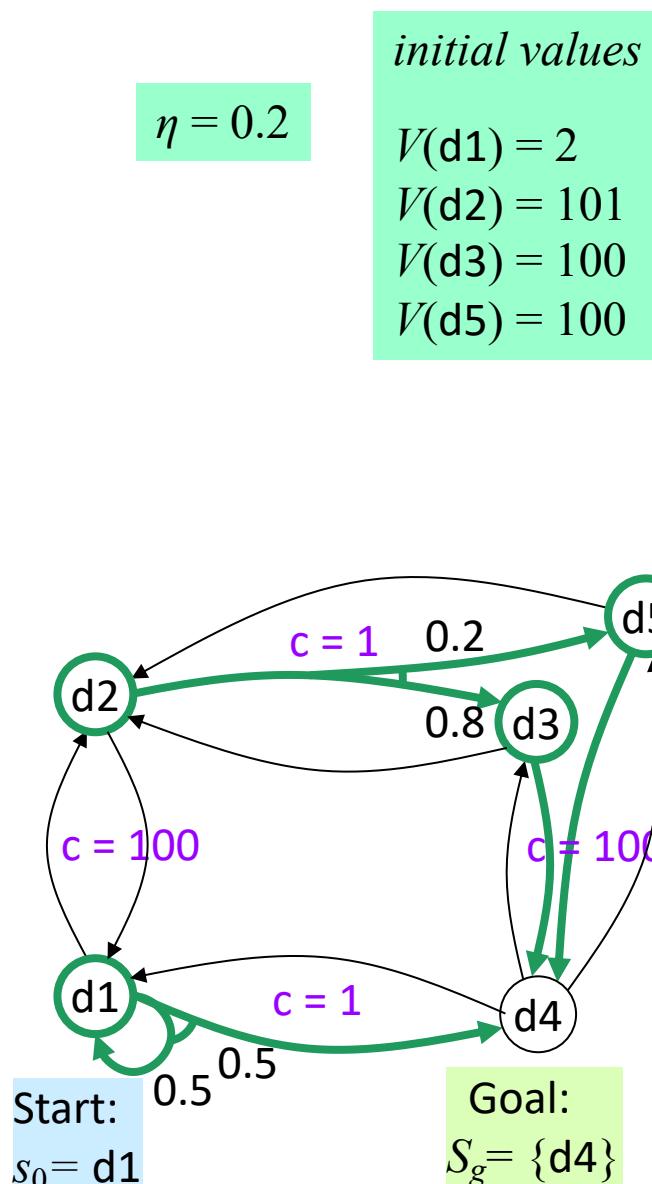
```

Bellman-Update(s)

```

global  $V, \pi$ 
 $v_{\text{old}} \leftarrow V(s)$ 
for every  $a \in \text{Applicable}(s)$  do
   $Q(s,a) \leftarrow \text{cost}(s,a) + \sum_{s' \in S} \Pr(s'|s,a) V(s')$ 
 $V(s) \leftarrow \min_{a \in \text{Applicable}(s)} Q(s,a)$ 
 $\pi(s) \leftarrow \operatorname{argmin}_{a \in \text{Applicable}(s)} Q(s,a)$ 
return  $|V(s) - v_{\text{old}}|$ 

```



$$Q(d1,m12) = 100 + 101 = 201$$

$$Q(d1,m14) = 1 + \frac{1}{2}(0) + \frac{1}{2}(2) = 2$$

$$V(d1) = 2; \pi(d1) = m14; 2 - 2 = 0$$

$$Q(d2,m21) = 100 + 2 = 102$$

$$Q(d2,m23) = 1 + .8(100) + .2(100) = 101$$

$$V(d2) = 101; \pi(d2) = m23; 101 - 101 = 0$$

$$Q(d3,m32) = 1 + 101 = 102$$

$$Q(d3,m34) = 100 + 0 = 100$$

$$V(d3) = 100; \pi(d3) = m34; 100 - 100 = 0$$

$$Q(d5,m52) = 1 + 101 = 102$$

$$Q(d5,m54) = 100 + 0 = 100$$

$$V(d5) = 100; \pi(d5) = m54; 100 - 100 = 0$$

$$r = \max(0, 0, 0, 0) = 0 < \eta$$

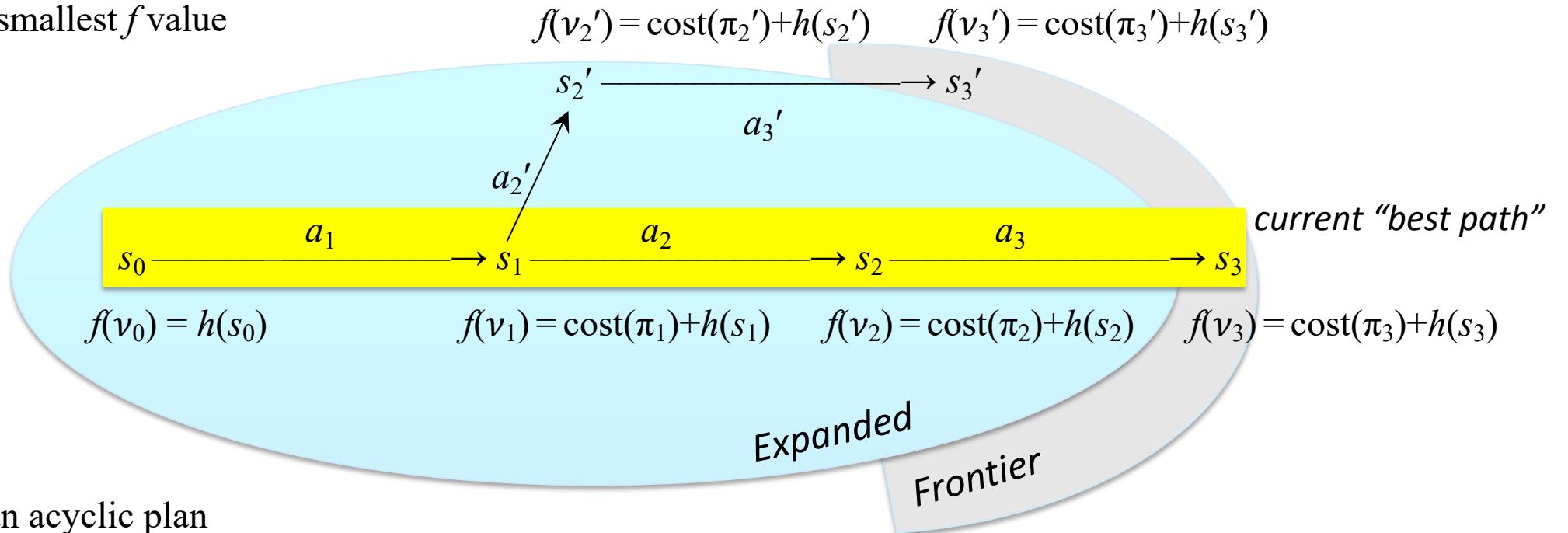
VI returns π

Discussion

- In each iteration of Policy Iteration
 - ▶ Compute V^π for current π
 - $|S|$ linear equations, $|S|$ unknowns
 - ▶ Use V^π to choose new π
 - $\text{argmin } Q$ value at each state
 - More work per iteration than value iteration
 - ▶ Needs to solve simultaneous equations
 - Usually converges in a smaller number of iterations
-
- In each iteration of Value Iteration
 - ▶ Compute new V : $\min Q$ value at each state
 - ▶ New V is a revised set of heuristic estimates
 - Doesn't depend on any π
 - ▶ Can let $\pi = \text{argmin } Q$ value at each state
 - But V isn't V^π for π or any other policy
 - ▶ Less work per iteration: doesn't need to solve a set of equations
 - ▶ Usually takes more iterations to converge
-
- At each iteration, both algorithms need to examine the entire state space
 - ▶ Number of iterations polynomial in $|S|$, but $|S|$ may be quite large
 - Next: use search techniques to avoid searching the entire space

Digression: A* Search

- All actions are deterministic
- Each node v_i is a pair (π_i, s_i)
- $f(v_i) = \text{cost of following } \pi_i \text{ from } s_0 \text{ to } s_i + h(s_i)$
- A* expands the frontier node that has the smallest f value

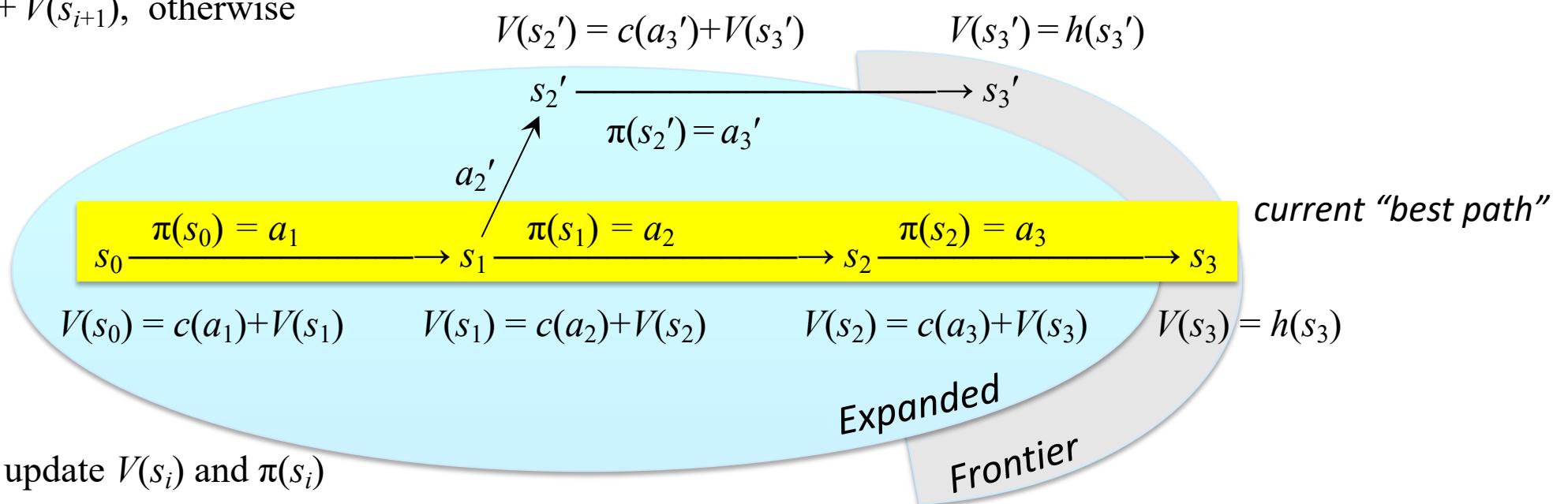


- A* creates an acyclic plan
- In deterministic planning domains, equivalent to an acyclic policy
 - ▶ Above, $\langle a_1, a_2, a_3 \rangle$ is equivalent to $\{(s_0, a_1), (s_1, a_2), (s_2, a_3)\}$
- We can rewrite A* to create an acyclic policy ...

Equivalent Approach Using Policies

- Global policy $\pi = \{(s_0, a_1), (s_1, a_2), (s_2, a_3), (s_3, a_4), (s_3', a_4')\}$
- $V(s_i) = \text{cost(following } \pi \text{ from } s_i \text{ to a frontier node}) + h(\text{the frontier node})$

$$= \begin{cases} h(s_i), & \text{if } s_i \text{ is a frontier node} \\ c(a_i) + V(s_{i+1}), & \text{otherwise} \end{cases}$$



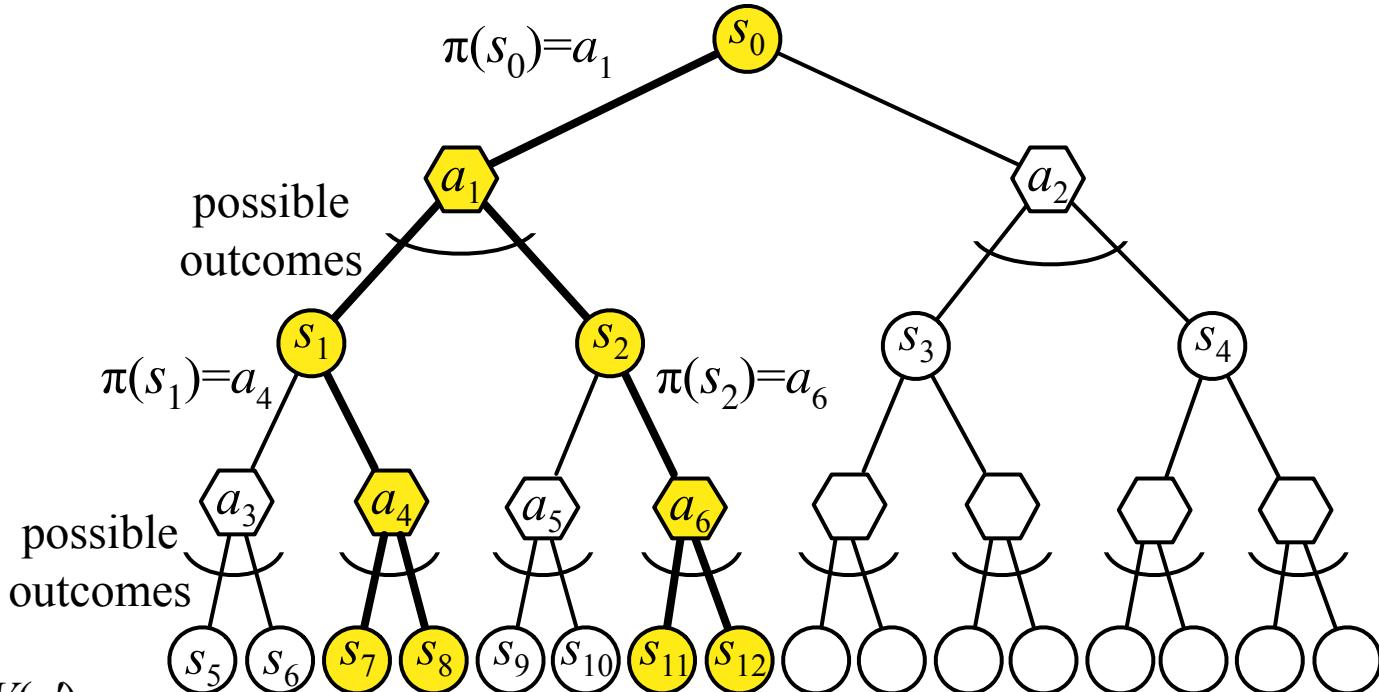
- Expand s_3
- For $i = 3, 2, 1, 0$, update $V(s_i)$ and $\pi(s_i)$
 - ▶ Updating:
 - $V(s) = \min_{a \in \text{Applicable}(s)} c(a) + V(\gamma(s, a))$
 - $\pi(s) = \operatorname{argmin}_{a \in \text{Applicable}(s)} c(a) + V(\gamma(s, a))$
 - ▶ E.g., if $c(a_2') + V(s_2') < c(a_2) + V(s_2)$ then $\pi(s_1) \leftarrow a_2'$

Poll: If $V(s)$ doesn't change, do we need to update its parent(s)?

A. yes
B. no

AO* (Basic Idea)

- An SSP can be represented as an AND/OR graph
 - OR nodes: choose an action
 - AND nodes: action's outcomes
- AO*: generalization of A* for *acyclic* SSPs
- $\text{leaves}(s_0, \pi) = \{s_7, s_8, s_{11}, s_{12}\}$
 - Expand one of them, e.g., s_{11}
- Going bottom-up, update V and π values for s_{11} and its ancestors
 - $V(s) = \min_{a \in \text{Applicable}(s)} c(a) + \sum_{s' \in \gamma(s,a)} P(s' | s, a) V(s')$
 - $\pi(s) = \operatorname{argmin}_{a \in \text{Applicable}(s)} c(a) + \sum_{s' \in \gamma(s,a)} P(s' | s, a) V(s')$
- e.g.,
 - $V(s_2) = \min(c(a_5)+V(s_9)+V(s_{10}), c(a_6)+V(s_{11})+V(s_{12}))$

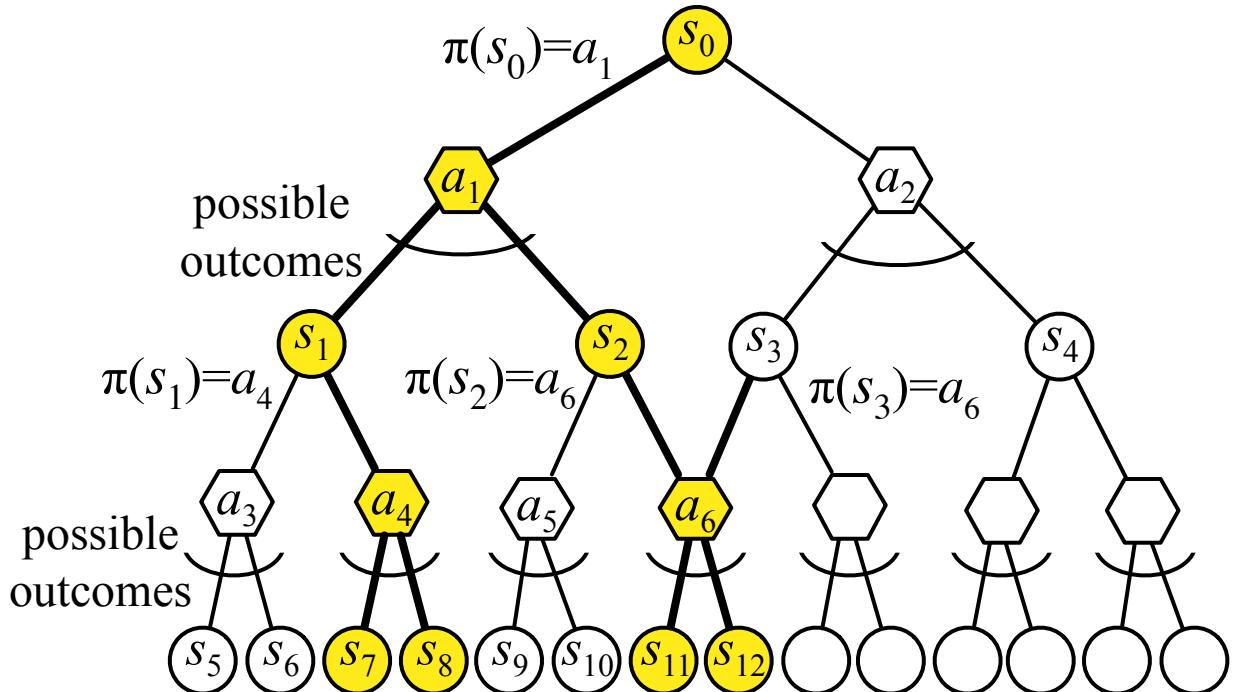


Poll: If $V(s)$ doesn't change, do we need to update its parent(s)?
 A. yes
 B. no

Poll: Will this work correctly if the acyclic AND/OR graph isn't a tree?
 A. yes
 B. no

AO* (Basic Idea)

- An SSP can be represented as an AND/OR graph
 - OR nodes: choose an action
 - AND nodes: action's outcomes
- AO*: generalization of A* for *acyclic* SSPs
- $\text{leaves}(s_0, \pi) = \{s_7, s_8, s_{11}, s_{12}\}$
 - Expand one of them, e.g., s_{11}
- Going bottom-up, update V and π values for s_{11} **and its ancestors**
 - $V(s) = \min_{a \in \text{Applicable}(s)} c(a) + \sum_{s' \in \gamma(s, a)} V(s')$
 - $\pi(s) = \operatorname{argmin}_{a \in \text{Applicable}(s)} c(a) + \sum_{s' \in \gamma(s, a)} V(s')$
- e.g.,
 - $V(s_2) = \min(c(a_5) + V(s_9) + V(s_{10}), c(a_6) + V(s_{11}) + V(s_{12}))$



Poll: If we change $V(s)$ and s has more than one parent, do we need to update all of them?

- yes
- no
- sometimes

AO*

$\text{AO}^*(\Sigma, s_0, S_g, V_0)$

global $\pi \leftarrow \emptyset$

global $\text{Envelope} \leftarrow \{s_0\}$

global $V; V(s_0) \leftarrow V_0(s_0)$

while $\text{Fringe} \neq \emptyset$ do

 select a state $s \in \text{Fringe}$

 for all $a \in \text{Applicable}(s)$ and $s' \in \gamma(s, a)$ do

 if $s' \notin \text{Envelope}$ then

 add s' to Envelope

$V(s') \leftarrow V_0(s')$

 AO-Update(s)

return π

AO-Update(s) // update V and π values of s and its ancestors

$Z \leftarrow \{s\}$ // states that need updating

while $Z \neq \emptyset$ do

 select $s \in Z$ such that $\hat{\gamma}(s, \pi(s)) \cap Z = \{s\}$

 remove s from Z

 Bellman-Update(s)

$Z \leftarrow Z \cup \{s' \in \text{Envelope} \mid s \in \gamma(s', \pi)\}$

a “lowest” state in Z

add s 's parents

acyclic

like *Expanded* \cup *Frontier* in A*

$\text{Fringe} = \text{Envelope} \setminus (\text{Dom}(\pi) \cup S_g)$

- like non-goal Frontier states in A*

• Fringe changes when
Bellman-Update changes π

Bellman-Update(s)

global V, π

$v_{\text{old}} \leftarrow V(s)$

for every $a \in \text{Applicable}(s)$ do

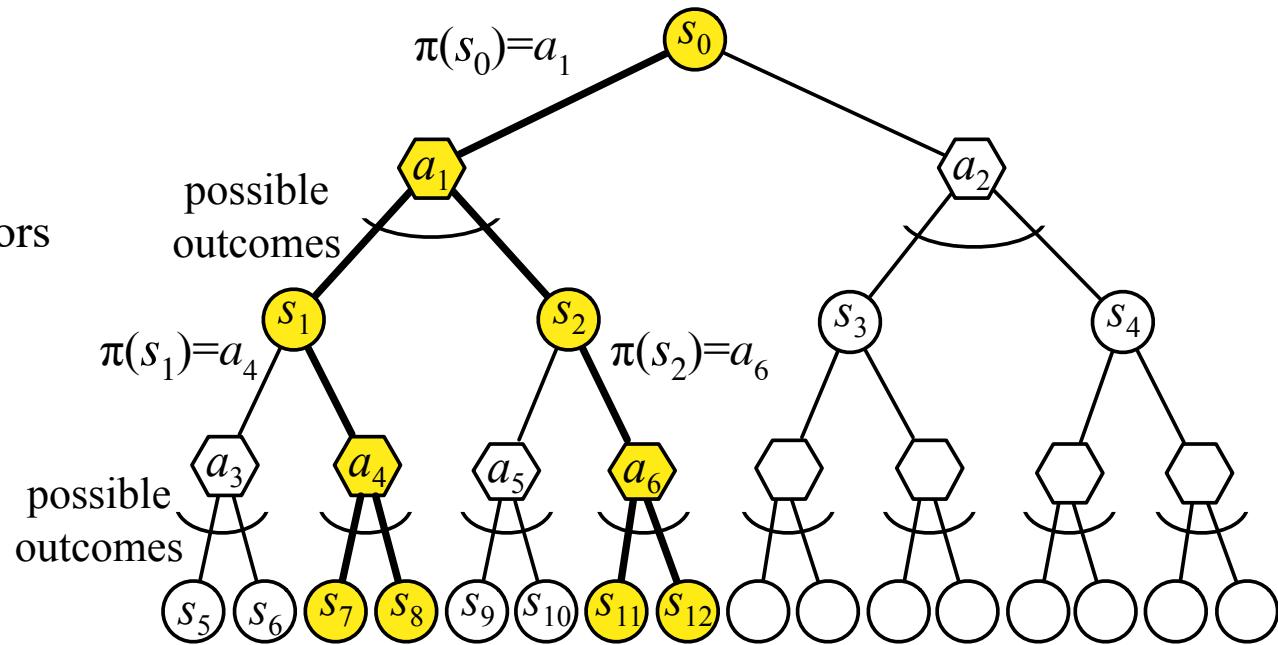
$Q(s, a) \leftarrow \text{cost}(s, a) + \sum_{s' \in S} \Pr(s'|s, a) V(s')$

$V(s) \leftarrow \min_{a \in \text{Applicable}(s)} Q(s, a)$

$\pi(s) \leftarrow \operatorname{argmin}_{a \in \text{Applicable}(s)} Q(s, a)$

return $|V(s) - v_{\text{old}}|$

not needed
this time



AO*

```

AO* ( $\Sigma, s_0, S_g, V_0$ )
    global  $\pi \leftarrow \emptyset$ 
    global  $Envelope \leftarrow \{s_0\}$ 
    global  $V; V(s_0) \leftarrow V_0(s_0)$ 
    while  $Fringe \neq \emptyset$  do
        select a state  $s \in Fringe$ 
        for all  $a \in Applicable(s)$  and  $s' \in \gamma(s,a)$  do
            if  $s' \notin Envelope$  then
                add  $s'$  to  $Envelope$ 
                 $V(s') \leftarrow V_0(s')$ 
        AO-Update( $s$ )
    return  $\pi$ 

```

```

AO-Update( $s$ ) // update  $V$  and  $\pi$  values of  $s$  and its ancestors
 $Z \leftarrow \{s\}$  // states that need updating
while  $Z \neq \emptyset$  do
    select  $s \in Z$  such that  $\hat{\gamma}(s, \pi(s)) \cap Z = \{s\}$ 
    remove  $s$  from  $Z$ 
    Bellman-Update( $s$ )
     $Z \leftarrow Z \cup \{s' \in Envelope \mid s \in \gamma(s', \pi)\}$ 

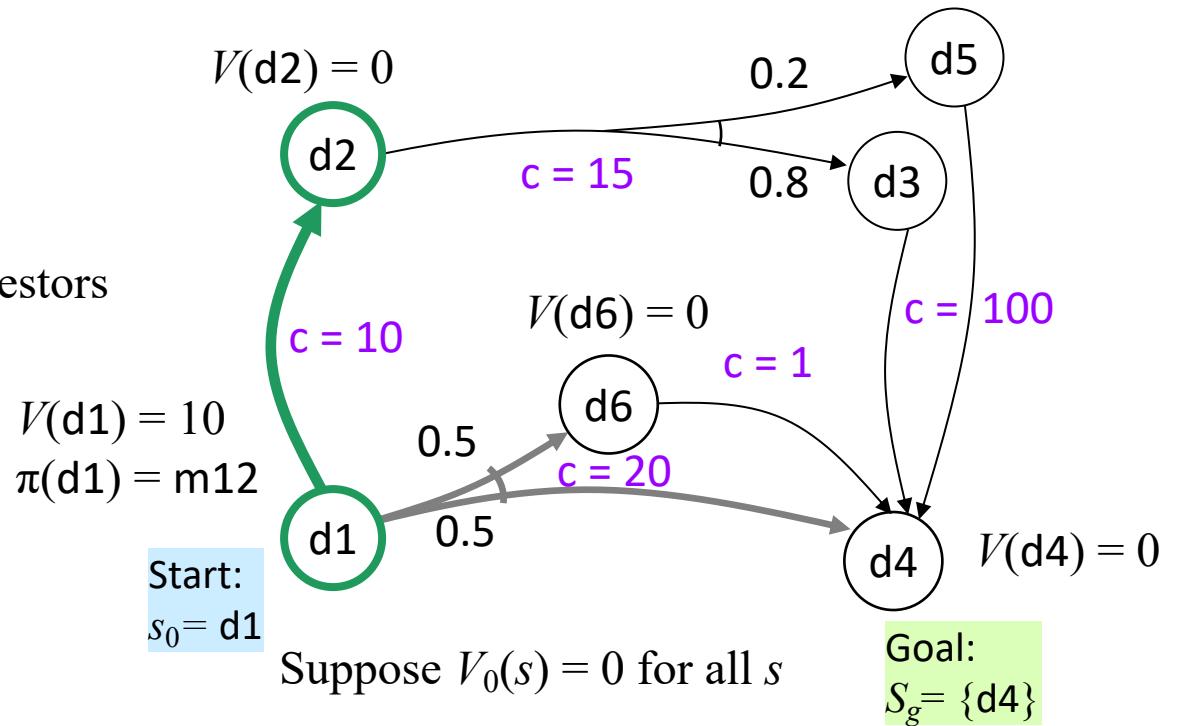
```

Bellman-Update(s)

```

    global  $V, \pi$ 
     $v_{old} \leftarrow V(s)$ 
    for every  $a \in Applicable(s)$  do
         $Q(s,a) \leftarrow \text{cost}(s,a) + \sum_{s' \in S} \Pr(s'|s,a) V(s')$ 
     $V(s) \leftarrow \min_{a \in Applicable(s)} Q(s,a)$ 
     $\pi(s) \leftarrow \operatorname{argmin}_{a \in Applicable(s)} Q(s,a)$ 
    return  $|V(s) - v_{old}|$ 

```



AO*

```

AO* ( $\Sigma, s_0, S_g, V_0$ )
    global  $\pi \leftarrow \emptyset$ 
    global  $Envelope \leftarrow \{s_0\}$ 
    global  $V; V(s_0) \leftarrow V_0(s_0)$ 
    while  $Fringe \neq \emptyset$  do
        select a state  $s \in Fringe$ 
        for all  $a \in Applicable(s)$  and  $s' \in \gamma(s,a)$  do
            if  $s' \notin Envelope$  then
                add  $s'$  to  $Envelope$ 
                 $V(s') \leftarrow V_0(s')$ 
        AO-Update( $s$ )
    return  $\pi$ 

```

```

AO-Update( $s$ ) // update  $V$  and  $\pi$  values of  $s$  and its ancestors
 $Z \leftarrow \{s\}$  // states that need updating
while  $Z \neq \emptyset$  do
    select  $s \in Z$  such that  $\hat{\gamma}(s, \pi(s)) \cap Z = \{s\}$ 
    remove  $s$  from  $Z$ 
    Bellman-Update( $s$ )
     $Z \leftarrow Z \cup \{s' \in Envelope \mid s \in \gamma(s', \pi)\}$ 

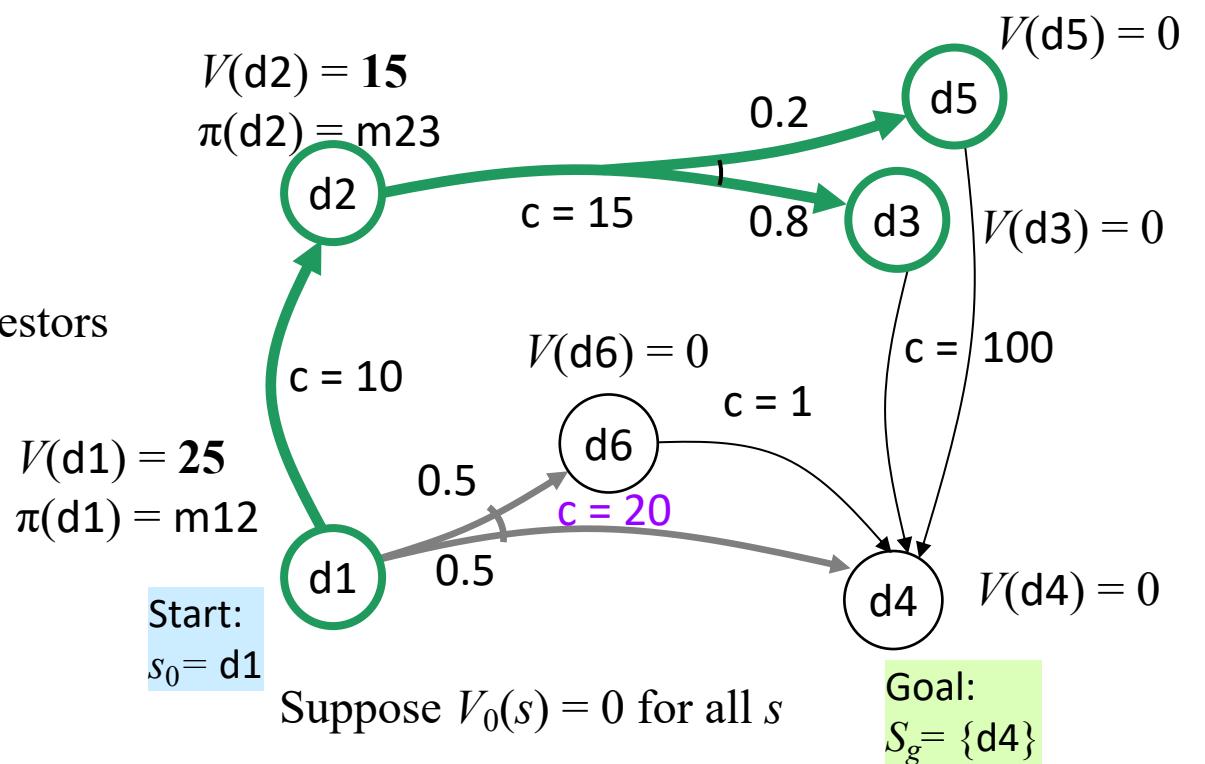
```

Bellman-Update(s)

```

    global  $V, \pi$ 
     $v_{old} \leftarrow V(s)$ 
    for every  $a \in Applicable(s)$  do
         $Q(s,a) \leftarrow \text{cost}(s,a) + \sum_{s' \in S} \Pr(s'|s,a) V(s')$ 
     $V(s) \leftarrow \min_{a \in Applicable(s)} Q(s,a)$ 
     $\pi(s) \leftarrow \operatorname{argmin}_{a \in Applicable(s)} Q(s,a)$ 
    return  $|V(s) - v_{old}|$ 

```



AO*

```

AO* ( $\Sigma, s_0, S_g, V_0$ )
    global  $\pi \leftarrow \emptyset$ 
    global  $Envelope \leftarrow \{s_0\}$ 
    global  $V; V(s_0) \leftarrow V_0(s_0)$ 
    while  $Fringe \neq \emptyset$  do
        select a state  $s \in Fringe$ 
        for all  $a \in Applicable(s)$  and  $s' \in \gamma(s,a)$  do
            if  $s' \notin Envelope$  then
                add  $s'$  to  $Envelope$ 
                 $V(s') \leftarrow V_0(s')$ 
            AO-Update( $s$ )
        return  $\pi$ 
    
```

```

AO-Update( $s$ ) // update  $V$  and  $\pi$  values of  $s$  and its ancestors
 $Z \leftarrow \{s\}$  // states that need updating
while  $Z \neq \emptyset$  do
    select  $s \in Z$  such that  $\hat{\gamma}(s, \pi(s)) \cap Z = \{s\}$ 
    remove  $s$  from  $Z$ 
    Bellman-Update( $s$ )
     $Z \leftarrow Z \cup \{s' \in Envelope \mid s \in \gamma(s', \pi)\}$ 

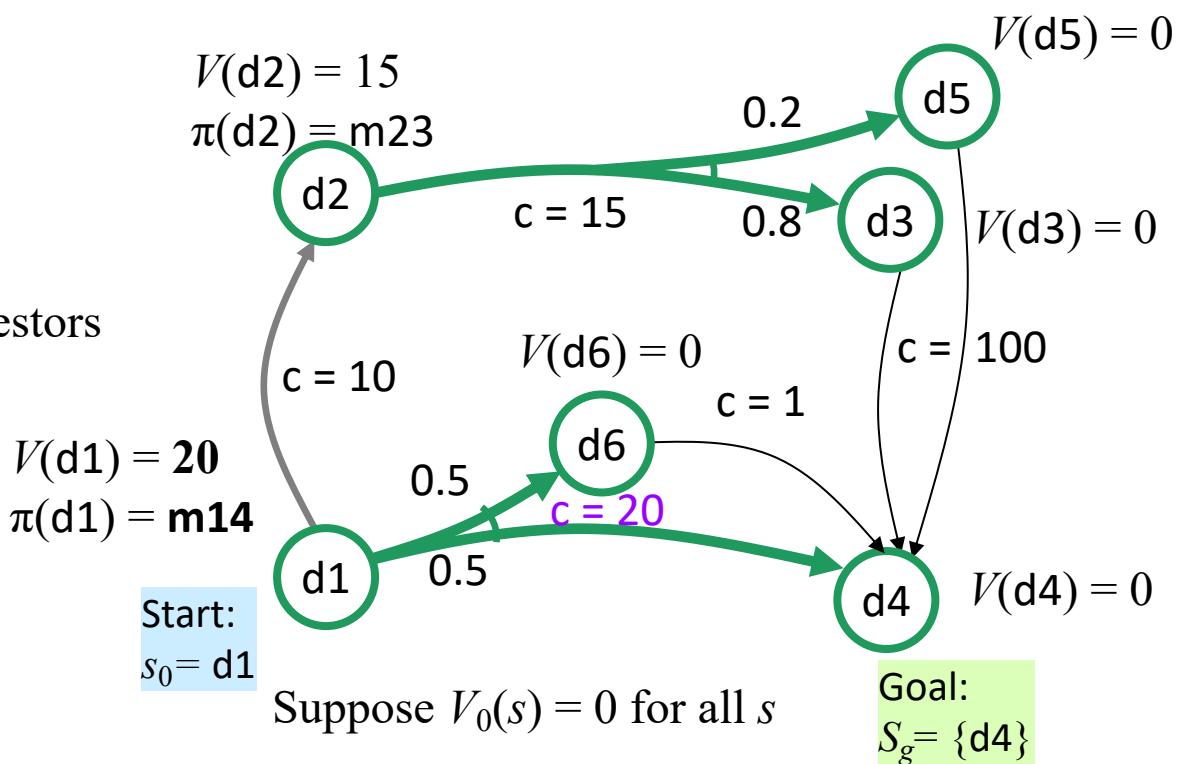
```

Bellman-Update(s)

```

    global  $V, \pi$ 
     $v_{old} \leftarrow V(s)$ 
    for every  $a \in Applicable(s)$  do
         $Q(s,a) \leftarrow \text{cost}(s,a) + \sum_{s' \in S} \Pr(s'|s,a) V(s')$ 
     $V(s) \leftarrow \min_{a \in Applicable(s)} Q(s,a)$ 
     $\pi(s) \leftarrow \operatorname{argmin}_{a \in Applicable(s)} Q(s,a)$ 
    return  $|V(s) - v_{old}|$ 

```



AO*

```

AO* ( $\Sigma, s_0, S_g, V_0$ )
    global  $\pi \leftarrow \emptyset$ 
    global  $Envelope \leftarrow \{s_0\}$ 
    global  $V; V(s_0) \leftarrow V_0(s_0)$ 
    while  $Fringe \neq \emptyset$  do
        select a state  $s \in Fringe$ 
        for all  $a \in Applicable(s)$  and  $s' \in \gamma(s,a)$  do
            if  $s' \notin Envelope$  then
                add  $s'$  to  $Envelope$ 
                 $V(s') \leftarrow V_0(s')$ 
            AO-Update( $s$ )
        return  $\pi$ 
    
```

```

AO-Update( $s$ ) // update  $V$  and  $\pi$  values of  $s$  and its ancestors
 $Z \leftarrow \{s\}$  // states that need updating
while  $Z \neq \emptyset$  do
    select  $s \in Z$  such that  $\hat{\gamma}(s, \pi(s)) \cap Z = \{s\}$ 
    remove  $s$  from  $Z$ 
    Bellman-Update( $s$ )
     $Z \leftarrow Z \cup \{s' \in Envelope \mid s \in \gamma(s', \pi)\}$ 

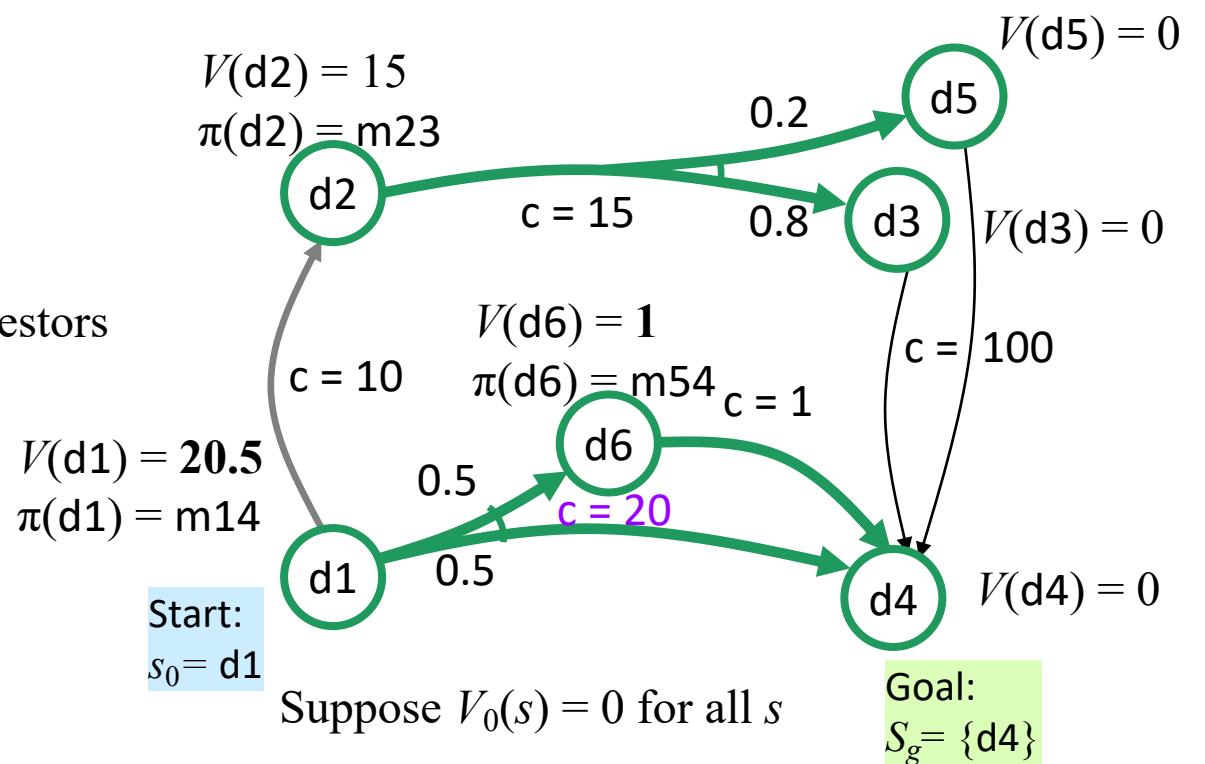
```

Bellman-Update(s)

```

    global  $V, \pi$ 
     $v_{old} \leftarrow V(s)$ 
    for every  $a \in Applicable(s)$  do
         $Q(s,a) \leftarrow \text{cost}(s,a) + \sum_{s' \in S} \Pr(s'|s,a) V(s')$ 
     $V(s) \leftarrow \min_{a \in Applicable(s)} Q(s,a)$ 
     $\pi(s) \leftarrow \operatorname{argmin}_{a \in Applicable(s)} Q(s,a)$ 
    return  $|V(s) - v_{old}|$ 

```



What to do about dead ends?

$\text{AO}^*(\Sigma, s_0, S_g, V_0)$

global $\pi \leftarrow \emptyset$

global $\text{Envelope} \leftarrow \{s_0\}$

global V ; $V(s_0) \leftarrow V_0(s_0)$

while $\text{Fringe} \neq \emptyset$ do

select a state $s \in \text{Fringe}$

for all $a \in \text{Applicable}(s)$ and $s' \in \gamma(s, a)$ do

if $s' \notin \text{Envelope}$ then

add s' to Envelope

$V(s') \leftarrow V_0(s')$

$\text{AO-Update}(s)$

return π

$\text{AO-Update}(s)$ // update V and π values of s and its ancestors

$Z \leftarrow \{s\}$ // states that need updating

while $Z \neq \emptyset$ do

select $s \in Z$ such that $\hat{\gamma}(s, \pi(s)) \cap Z = \{s\}$

remove s from Z

$\text{Bellman-Update}(s)$

$Z \leftarrow Z \cup \{s' \in \text{Envelope} \mid s \in \gamma(s', \pi)\}$

$\text{Bellman-Update}(s)$

global V, π

$v_{\text{old}} \leftarrow V(s)$

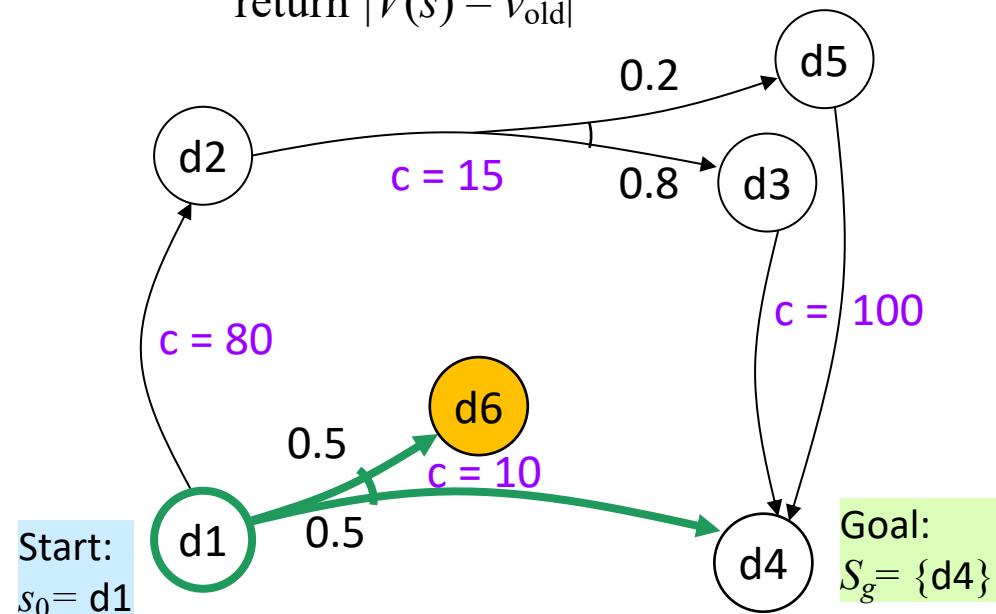
for every $a \in \text{Applicable}(s)$ do

$Q(s, a) \leftarrow \text{cost}(s, a) + \sum_{s' \in S} \Pr(s'|s, a) V(s')$

$V(s) \leftarrow \min_{a \in \text{Applicable}(s)} Q(s, a)$

$\pi(s) \leftarrow \operatorname{argmin}_{a \in \text{Applicable}(s)} Q(s, a)$

return $|V(s) - v_{\text{old}}|$

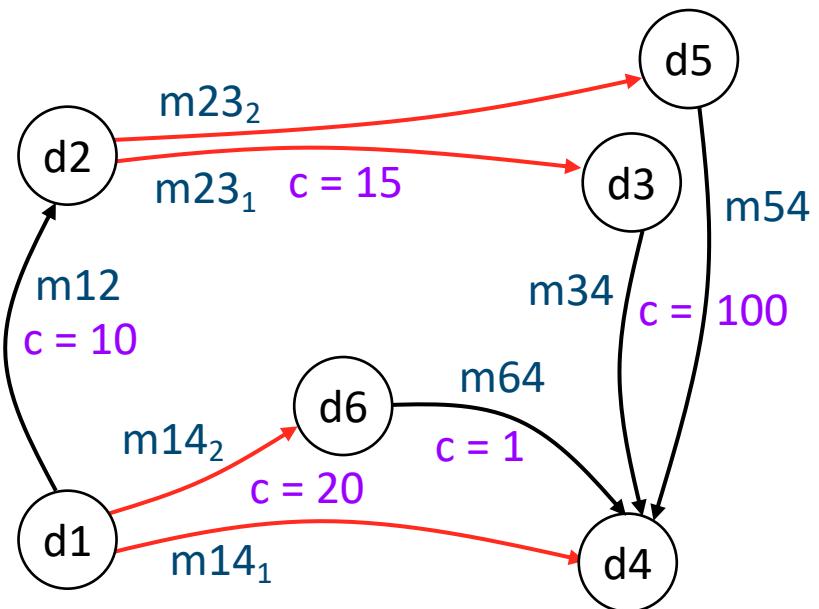
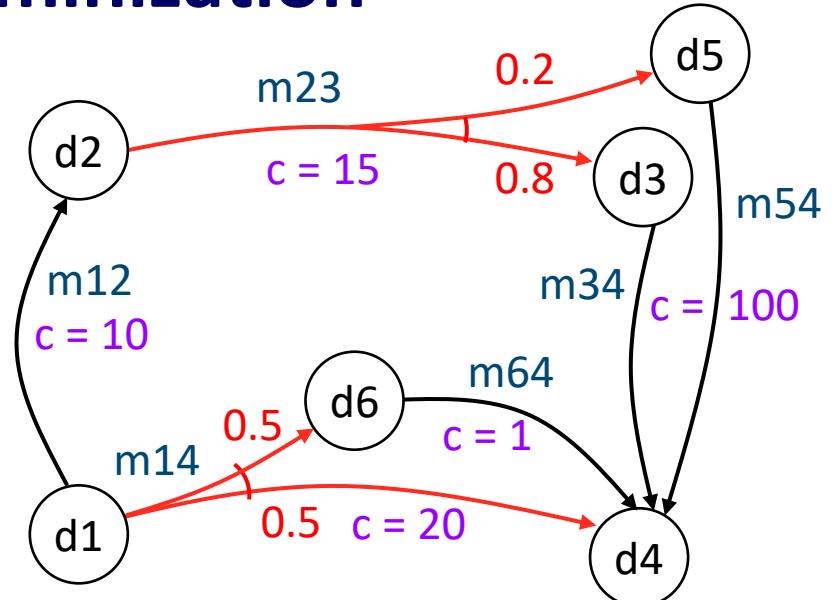


If action a in state s goes to a dead end, Guided-Find-Safe-Solution (Chap. 5) makes a inapplicable in s . AO^* can be modified to do this.

Heuristics through Determinization

What to use for $V_0(s)$?

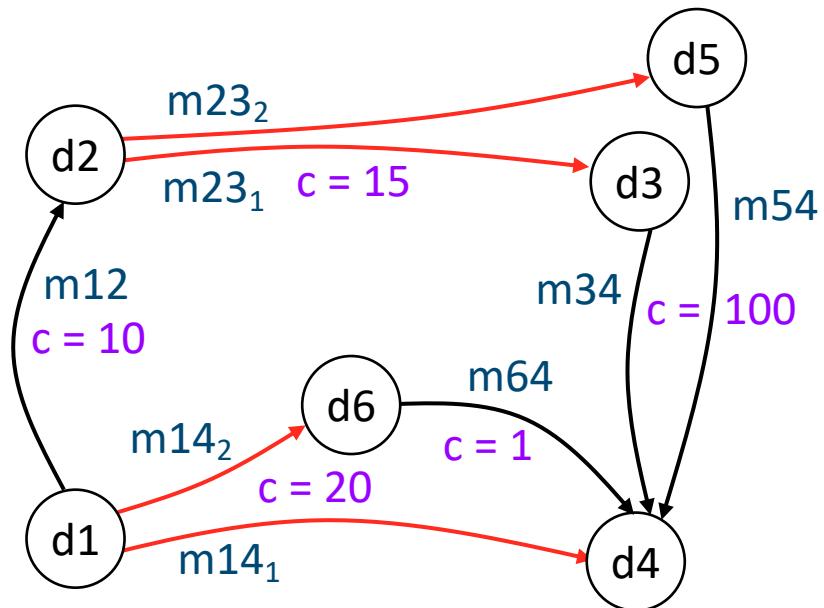
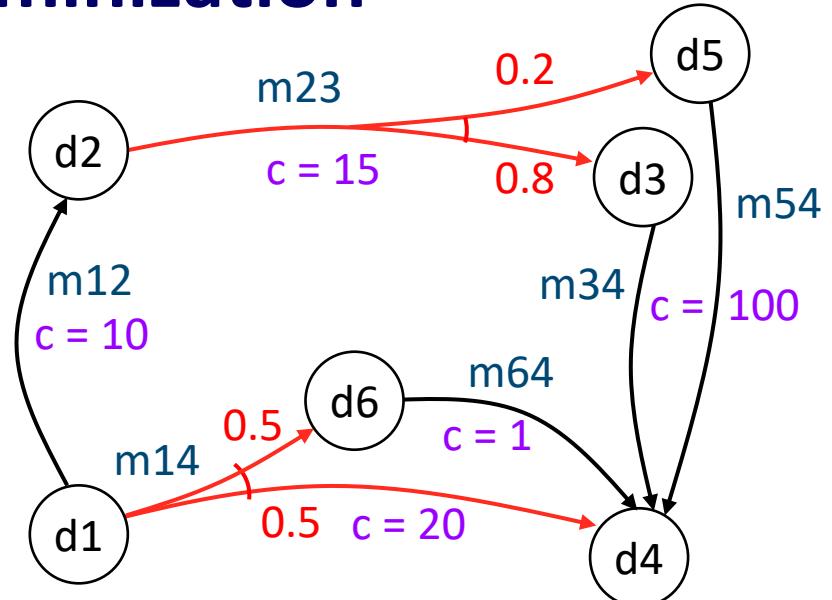
- One possibility: use a heuristic function for the *determinized* domain Σ_d
- For each action a of Σ
 - ▶ If a has k possible outcomes, then
 - Σ_d contains k deterministic actions a_1, \dots, a_k
 - one for each of a 's outcome
- Let h be a classical heuristic function for Σ_d
- If h is admissible for Σ_d then also admissible for Σ
- Why:
 - ▶ Let π be any optimal solution for Σ
 - ▶ Let p be any acyclic execution of π
 - p is a plan in Σ_d
 - $h(s_0) \leq \text{cost}(p) \leq \text{cost}(\pi)$



Heuristics through Determinization

What to use for $V_0(s)$?

- Another possibility: call a classical planner on the determinized problem (Σ_d, s, g)
 - Get plan $p = \langle a_1, a_2, \dots, a_n \rangle$
 - Return $V_0(s) = \text{cost}(p)$
- **Theorem.** If the classical planner always returns optimal plans, then V_0 is admissible
- Outline of proof:
 - ▶ Let π^* be an optimal solution for (Σ, s, g)
 - ▶ Then $V^*(s) = E(\text{cost}(\pi^*))$
= weighted avg. of all executions of π
 - ▶ Let p^* be the least costly execution of π
 - p^* is a solution for (Σ_d, s, g) , so the classical planner returns a plan p of $\leq \text{cost}$
 - Thus $V_0(s) = \text{cost}(p) \leq \text{cost}(p^*) \leq E(\text{cost}(\pi))$



LAO*

```

LAO* ( $\Sigma, s_0, S_g, V_0$ )
global  $\pi \leftarrow \emptyset$ 
global Envelope  $\leftarrow \{s_0\}$ 
global  $V$ ;  $V(s_0) \leftarrow V_0(s_0)$ 
while  $Fringe \neq \emptyset$  do
    select a state  $s \in Fringe$ 
    for all  $a \in Applicable(s)$  and  $s' \in \gamma(s,a)$  do
        if  $s' \notin Envelope$  then
            add  $s'$  to Envelope
             $V(s') \leftarrow V_0(s')$ 
        LAO-Update( $s$ )
    return  $\pi$ 

LAO-Update( $s$ )
 $Z \leftarrow \{s\} \cup \{s' \in Envelope \mid s \in \hat{\gamma}(s',\pi)\}$ 
loop until  $r \leq \eta$  or new states added to  $Fringe$ 
 $r \leftarrow \max\{\text{Bellman-Update}(s) \mid s \in Z\}$ 

```

may be cyclic

All π -ancestors of s in Envelope

Value iteration, restricted to Z

Same as AO* except here

Bellman-Update(s)

global V, π

$v_{\text{old}} \leftarrow V(s)$

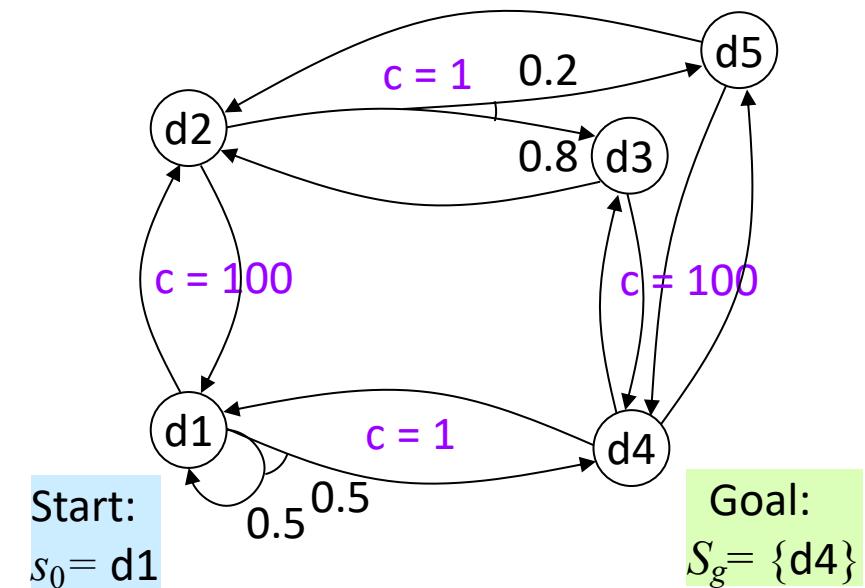
for every $a \in \text{Applicable}(s)$ do

$Q(s,a) \leftarrow \text{cost}(s,a) + \sum_{s' \in S} \Pr(s'|s,a) V(s')$

$V(s) \leftarrow \min_{a \in \text{Applicable}(s)} Q(s,a)$

$\pi(s) \leftarrow \operatorname{argmin}_{a \in \text{Applicable}(s)} Q(s,a)$

return $|V(s) - v_{\text{old}}|$



Example: $V_0(s) = 0$ for all s

LAO* (Σ, s_0, S_g, V_0)

global $\pi \leftarrow \emptyset$

global $Envelope \leftarrow \{s_0\}$

global V ; $V(s_0) \leftarrow V_0(s_0)$

while $Fringe \neq \emptyset$ do

select a state $s \in Fringe$

for all $a \in Applicable(s)$ and $s' \in \gamma(s,a)$ do

if $s' \in \gamma(s,a) \notin Envelope$ then

add s' to $Envelope$

$V(s') \leftarrow V_0(s')$

LAO-Update(s)

return π

LAO-Update(s)

$Z \leftarrow \{s\} \cup \{s' \in Envelope \mid s \in \hat{\gamma}(s',\pi)\}$

loop until $r \leq \eta$ or new states added to $Fringe$

$r \leftarrow \max \{Bellman-Update(s) \mid s \in Z\}$

Bellman-Update(s)

global V, π ; $v_{old} \leftarrow V(s)$

for every $a \in Applicable(s)$ do

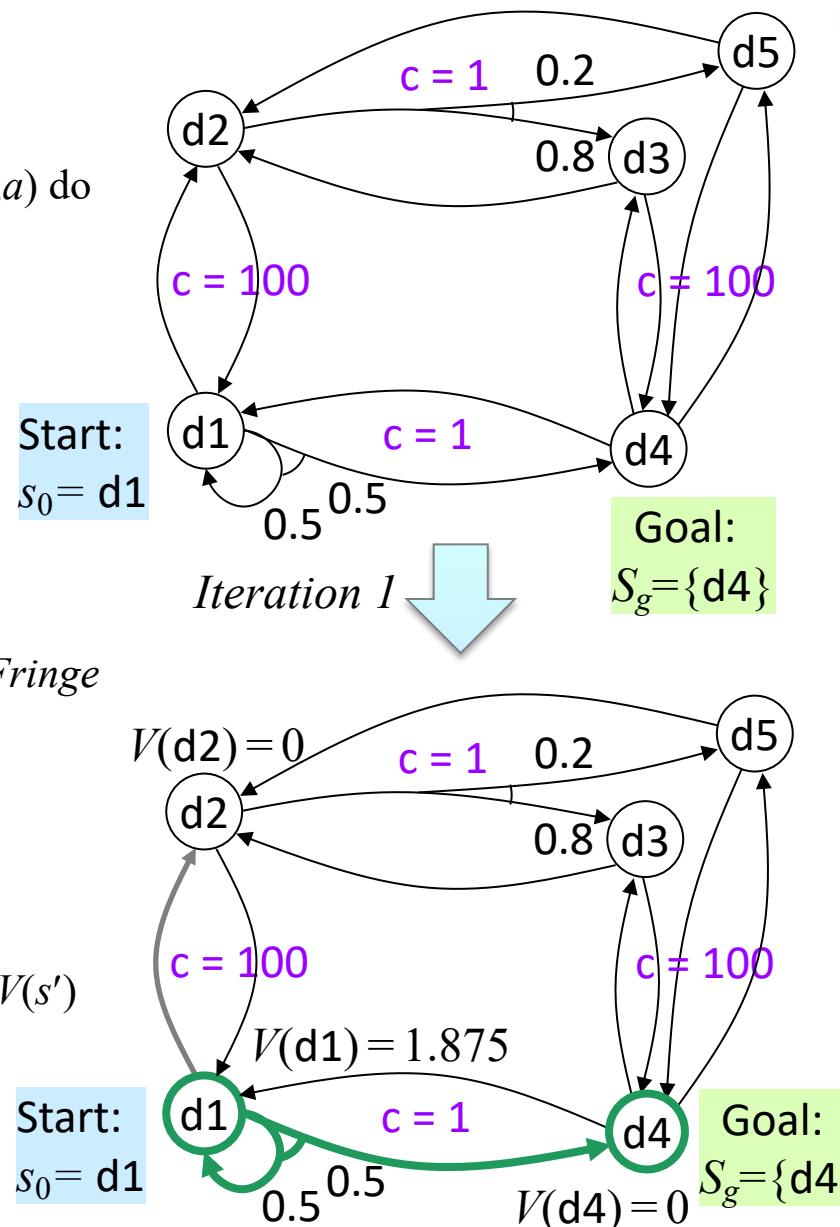
$Q(s,a) \leftarrow cost(s,a) + \sum_{s' \in S} \Pr(s'|s,a) V(s')$

$V(s) \leftarrow \min_{a \in Applicable(s)} Q(s,a)$

$\pi(s) \leftarrow \operatorname{argmin}_{a \in Applicable(s)} Q(s,a)$

return $|V(s) - v_{old}|$

Example 1



$$\eta = 0.2$$

$$V_0(s) = 0 \text{ for all } s$$

$$Envelope = \{d1\}$$

$$\pi = \emptyset$$

LAO* iteration 1:

select $s = d1$

Applicable($d1$) = {m12, m14}

add $d2$ to $Envelope$; $V(d2) \leftarrow 0$

add $d4$ to $Envelope$; $V(d4) \leftarrow 0$

Call LAO-Update($d1$)

π is empty, so $Z = \{d1\}$

LAO-Update iteration 1:

Call Bellman-update($d1$):

$$Q(d1, m12) = 100 + 0 = 100$$

$$Q(d1, m14) = 1 + (\frac{1}{2}(0) + \frac{1}{2}(0)) = 1$$

$$V(d1) = 1; \pi(d1) = m14$$

$$r = V(d1) - 0 = 1$$

Keep iterating until $r \leq 0.2$

$$V(d1) = 1.875; r = 0.0625$$

LAO* iteration 2:

$Fringe = \emptyset$, so return $\pi = \{(d1, m14)\}$

LAO* (Σ, s_0, S_g, V_0)

global $\pi \leftarrow \emptyset$

global $Envelope \leftarrow \{s_0\}$

global V ; $V(s_0) \leftarrow V_0(s_0)$

while $Fringe \neq \emptyset$ do

select a state $s \in Fringe$

for all $a \in Applicable(s)$ and $s' \in \gamma(s,a)$ do

if $s' \in \gamma(s,a) \notin Envelope$ then

add s' to $Envelope$

$V(s') \leftarrow V_0(s')$

LAO-Update(s)

return π

LAO-Update(s)

$Z \leftarrow \{s\} \cup \{s' \in Envelope \mid s \in \hat{\gamma}(s',\pi)\}$

loop until $r \leq \eta$ or new states added to $Fringe$

$r \leftarrow \max \{Bellman-Update(s) \mid s \in Z\}$

Bellman-Update(s)

global V, π ; $v_{old} \leftarrow V(s)$

for every $a \in Applicable(s)$ do

$Q(s,a) \leftarrow cost(s,a) + \sum_{s' \in S} \Pr(s'|s,a) V(s')$

$V(s) \leftarrow \min_{a \in Applicable(s)} Q(s,a)$

$\pi(s) \leftarrow \operatorname{argmin}_{a \in Applicable(s)} Q(s,a)$

return $|V(s) - v_{old}|$

Example 2

$\eta=0.2$; $V_0(s) = 0$ for all s

Envelope = {d1}; $\pi = \emptyset$

LAO* iteration 1:

select $s = d1$; Applicable(d1) = {m12, m14}

$V(d2) \leftarrow 0$; add d2 to Envelope

$V(d4) \leftarrow 0$; add d4 to Envelope

Call LAO-Update(d1)

π is empty, so $Z = \{d1\}$

LAO-Update iteration 1:

Call Bellman-update(d1):

$$Q(d1, m12) = 100 + 0 = 100$$

$$Q(d1, m14) = 80 + (\frac{1}{2}(0) + \frac{1}{2}(0)) = 80$$

$$V(d1) = 80; \pi(d1) = m14$$

$$r = V(d1) - 0 = 80$$

LAO-Update iteration 2:

Call Bellman-update(d1):

$$Q(d1, m12) = 100 + 0 = 100$$

$$Q(d1, m14) = 80 + (\frac{1}{2}(80) + \frac{1}{2}(0)) = 120$$

$$V(d1) = 100; \pi(d1) = m12$$

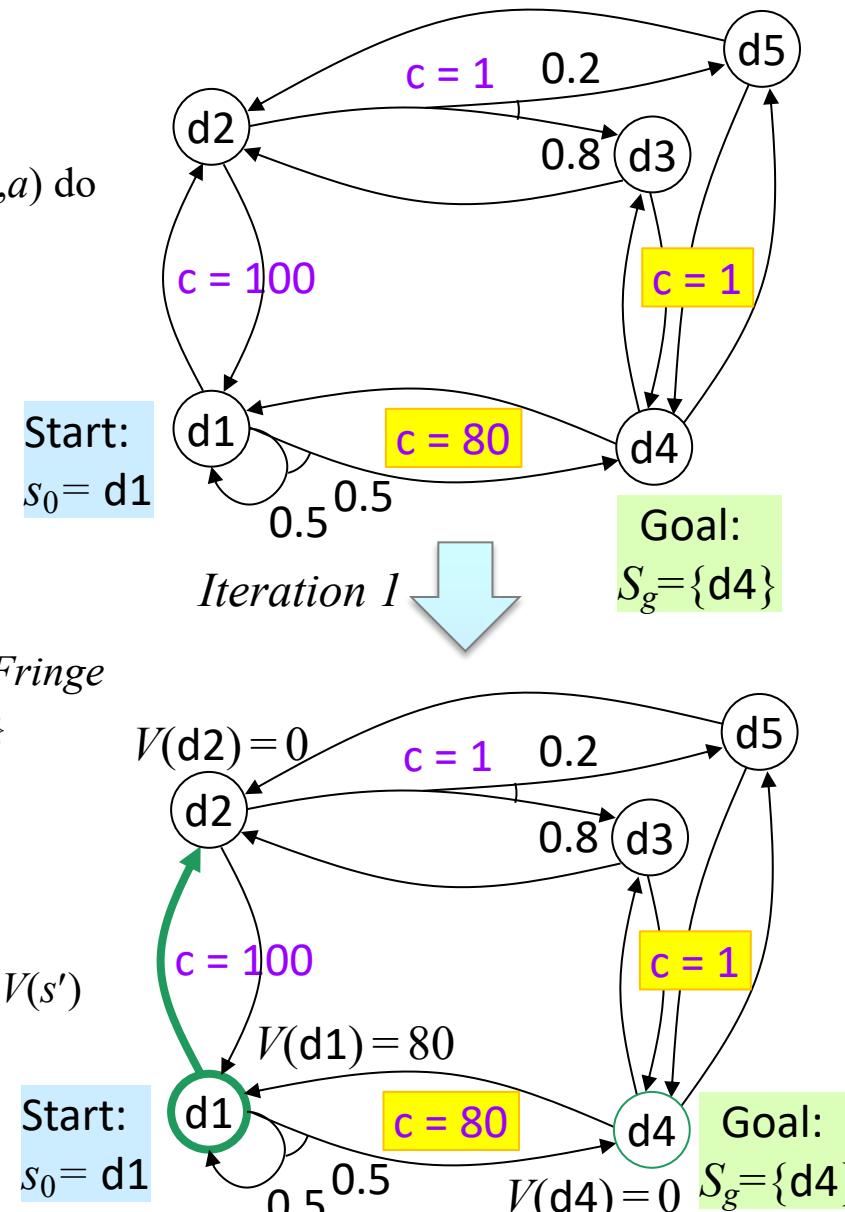
$$r = V(d1) - 80 = 20$$

new state d2 in Fringe, so return

LAO-Update returns

After more iterations, LAO* eventually returns

$$\pi = \{(d1, m12), (d2, m23), (d3, m34), (d5, m54)\}$$



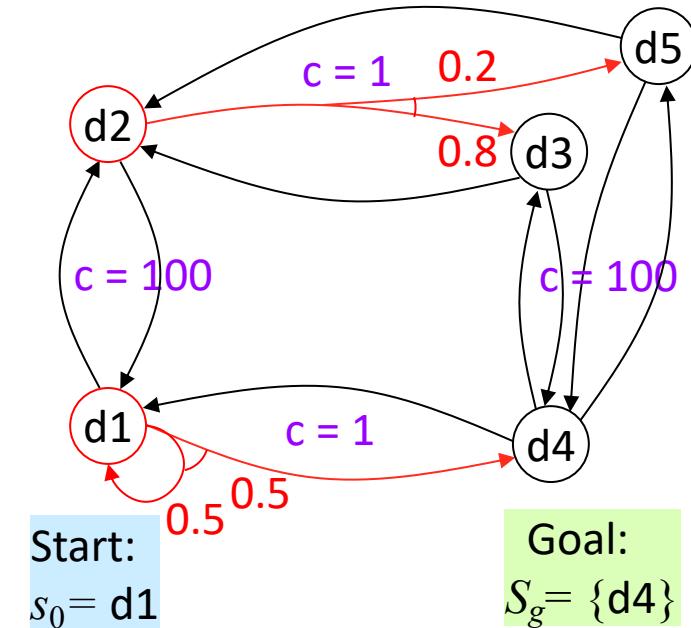
Skipping Ahead

- Skipping ILAO*, HDP, LD F S_a, LRTDP, SLATE
 - ▶ I'll come back to these if there's time

Planning and Acting

Run-Lookahead(Σ, s_0, S_g)

```
 $s \leftarrow s_0$ 
while  $s \notin S_g$  and Applicable( $s$ )  $\neq \emptyset$  do
   $a \leftarrow \text{Lookahead}(s, \theta)$ 
  perform action  $a$ 
   $s \leftarrow \text{observe resulting state}$ 
```



- Like Run-Lookahead from Chapter 2
- Differences:
 - Explicit starting state s_0
 - Not necessary, could observe s_0 instead
 - Doesn't abstract s (to simplify the presentation)
 - Lookahead returns an action instead of a plan

- What to use for Lookahead?
 - AO*, LAO*, ...
 - Modify to search part of the space
 - Classical planner on determinized domain
 - next page
 - Stochastic sampling algorithms

Planning and Acting

From Chapter 5

$\text{FS-Replan}(\Sigma, s, S_g)$

$\pi_d \leftarrow \emptyset$

while $s \notin S_g$ and $\text{Applicable}(s) \neq \emptyset$ do

if $\pi_d(s)$ is undefined then do

$\pi_d \leftarrow \text{Plan2policy}(\text{Forward-search } (\Sigma_d, s, S_g))$

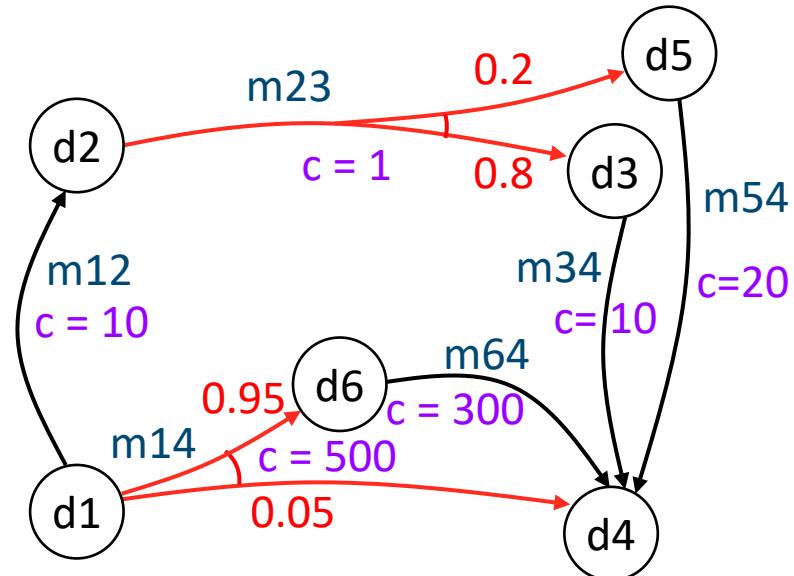
if $\pi_d = \text{failure}$ then return failure

perform action $\pi_d(s)$

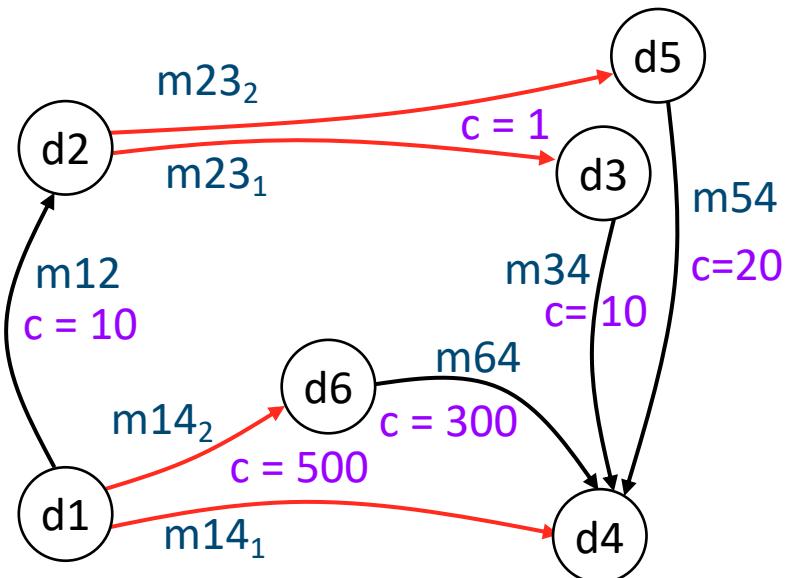
$s \leftarrow \text{observe resulting state}$

From Chapter 2

Determinized domain



- Generalization of a well-known planner called FF-Replan
- Like Run-Lazy-lookahead from Chapter 2
 - ▶ Lookahead = classical planner on determinized domain
- Example:
 - ▶ Forward-search returns $\langle m12, m23_1, m34 \rangle$
 - ▶ Plan2policy returns $\pi_d = \langle (d1, m12), (d2, m23), (d3, m34) \rangle$
 - ▶ If $m23$ goes to $d5$, then call Forward-search again



Planning and Acting

From Chapter 5

$\text{FS-Replan}(\Sigma, s, S_g)$

$$\pi_d \leftarrow \emptyset$$

while $s \notin S_g$ and $\text{Applicable}(s) \neq \emptyset$ do

if $\pi_d(s)$ is undefined then do

$\pi_d \leftarrow \text{Plan2policy}(\text{Forward-search } (\Sigma_d, s, S_g))$

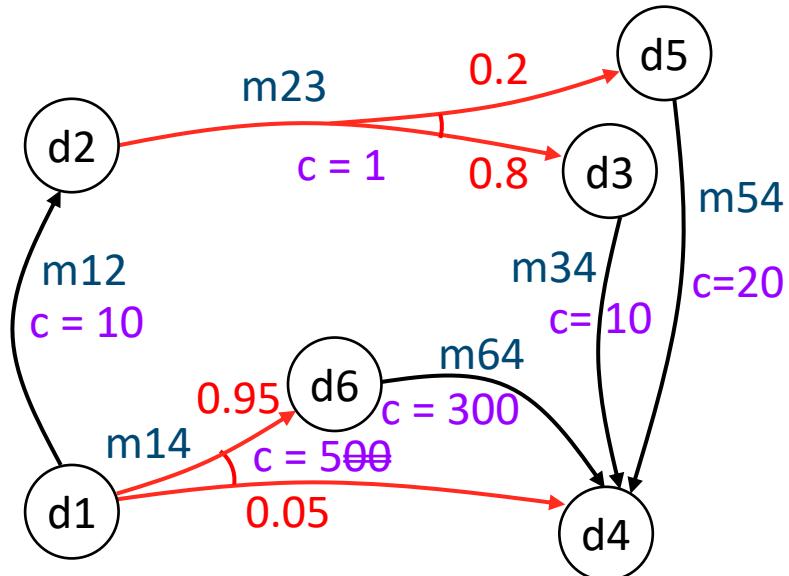
if $\pi_d = \text{failure}$ then return failure

perform action $\pi_d(s)$

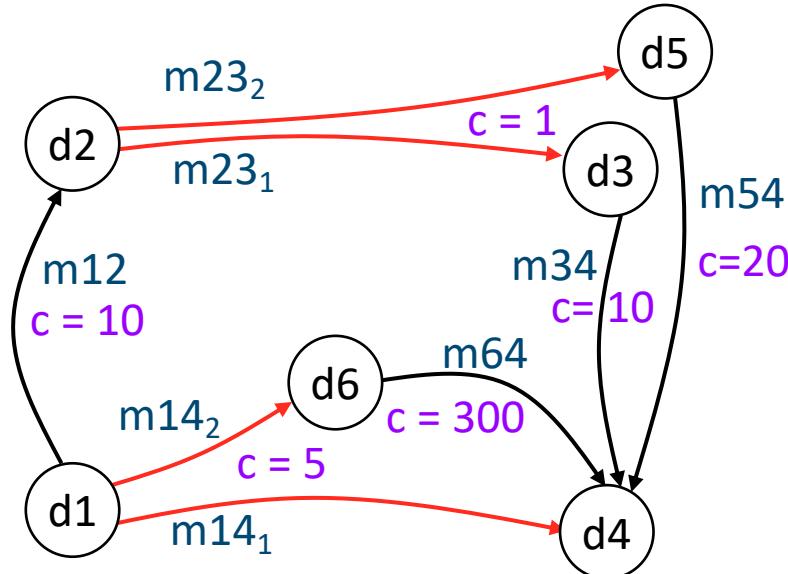
$s \leftarrow \text{observe resulting state}$

From Chapter 2

Determinized domain

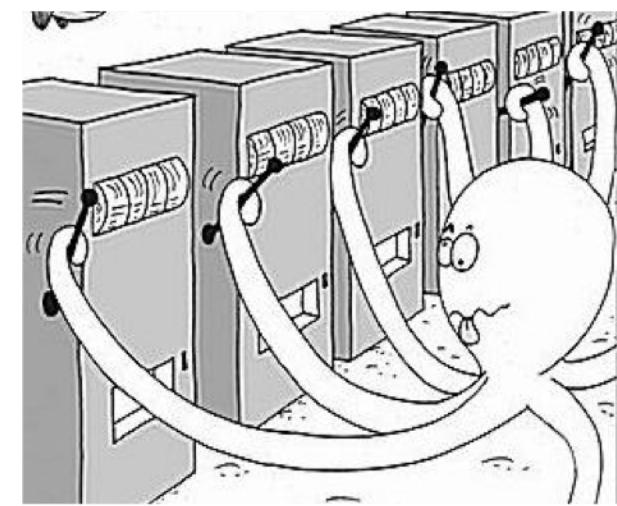


- Problem: classical planner may choose a plan that depends on low-probability outcome
- Example: Forward-search returns $\langle m14_1 \rangle$
 - ▶ $m14$ very likely to go to $d6 \Rightarrow$ large cost
- RFF algorithm (see book) attempts to alleviate this



Digression: Monte Carlo rollouts

- Multi-arm bandit problem
 - ▶ Statistical model of sequential experiments
 - ▶ Name derived from *one-armed bandit* (slot machine)
- Multiple actions a_1, a_2, \dots, a_n
 - ▶ Each a_i provides a reward from an unknown probability distribution p_i
 - ▶ Assume each p_i is *stationary*
 - Same every time, regardless of history
 - ▶ Objective: maximize expected utility of a sequence of actions
- Exploitation vs exploration dilemma:
 - ▶ **Exploitation:** choose an action that has given you high rewards in the past
 - ▶ **Exploration:** choose an action that's less familiar, in hopes that it might produce a higher reward



UCB (Upper Confidence Bound) Algorithm

- Assume all rewards are between 0 and 1
 - ▶ If they aren't, normalize them
- For each action a , let
 - ▶ $r(a)$ = average reward you've gotten from a
 - ▶ $n(a)$ = number of times you've tried a
 - ▶ $n_t = \sum_a n(a)$
 - ▶ $Q(a) = r(a) + \sqrt{2(\ln n_t)/n(a)}$

UCB:

if there are any untried actions:

$\tilde{a} \leftarrow$ any untried action

else:

$\tilde{a} \leftarrow \operatorname{argmax}_a Q(a)$

perform \tilde{a}

update $r(\tilde{a}), n(\tilde{a}), n_t, Q(\tilde{a})$

- Theorem (given some assumptions):
As the number of calls to UCB $\rightarrow \infty$,
UCB's choice at each call \rightarrow optimal choice

Actions:	a1	a2	a3	
No. of tries:	$n(a1) = 5$	$n(a2) = 3$	$n(a3) = 2$	$t = 10$
Rewards:	$r(a1) = 0.4$	0.3333	0	
Q values:	$Q(a1) = 1.35971$	1.5723	1.5174	
Payoffs:				
1st	0			
2nd		0		
3rd			0	
4th	1			
5th	0			
6th		0		
7th			0	
8th	1			
9th	0			
10th		1		

UCT Algorithm

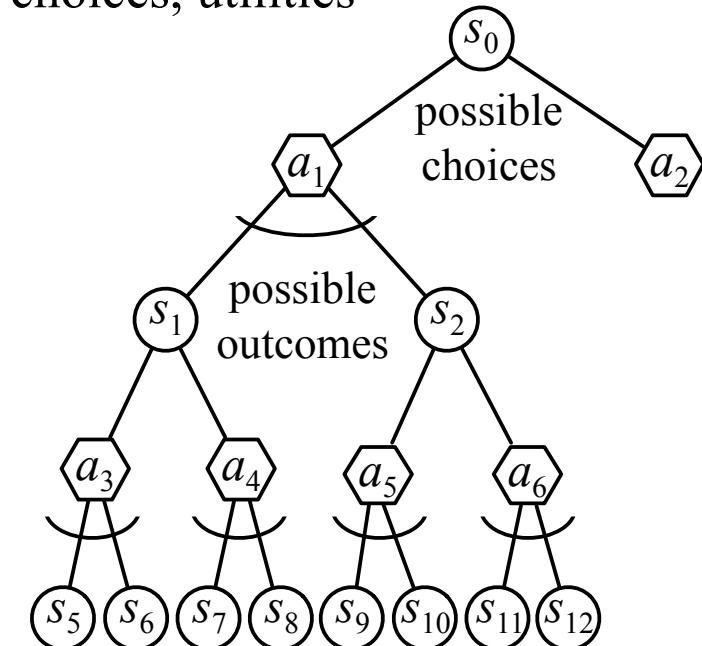
```

max search depth
UCT(s, h)
if  $s \in S_g$  then return 0
if  $h = 0$  then return  $V_0(s)$ 
if  $s \notin Envelope$  then do
    add  $s$  to  $Envelope$ 
     $n(s) \leftarrow 0$  like  $n_t$  on previous slide
    for all  $a \in Applicable(s)$  do
         $Q(s, a) \leftarrow 0; n(s, a) \leftarrow 0$  like  $Q(a), n(a)$  on previous slide
     $Untried \leftarrow \{a \in Applicable(s) \mid n(s, a) = 0\}$ 
    if  $Untried \neq \emptyset$  then  $\tilde{a} \leftarrow Choose(Untried)$ 
    else  $\tilde{a} \leftarrow argmin_{a \in Applicable(s)} \{Q(s, a) - C \times [\log(n(s)) / n(s, a)]^{1/2}\}$ 
     $s' \leftarrow Sample(\Sigma, s, \tilde{a})$ 
     $cost-rollout \leftarrow cost(s, \tilde{a}) + UCT(s', h - 1)$ 
     $Q(s, \tilde{a}) \leftarrow [n(s, \tilde{a}) \times Q(s, \tilde{a}) + cost-rollout] / (1 + n(s, \tilde{a}))$ 
     $n(s) \leftarrow n(s) + 1$ 
     $n(s, \tilde{a}) \leftarrow n(s, \tilde{a}) + 1$ 
return  $cost-rollout$ 

```

ln, I think
 $\sqrt{2}$ on previous slide
choose s' with probability $Pr(s' \mid s, \tilde{a})$

- $UCT(s, h)$: recursive UCB computation on an SSP
 - ▶ Adapted for minimization rather than maximization
- *Monte Carlo rollout*:
 - ▶ At s , choose action \tilde{a} using UCB computation
 - Perform \tilde{a} , get state s'
 - Do the same thing recursively at s'
 - Continue until reaching a goal, dead end, or depth h
 - ▶ At each state visited, keep statistics on choices, utilities



Using UCT Offline

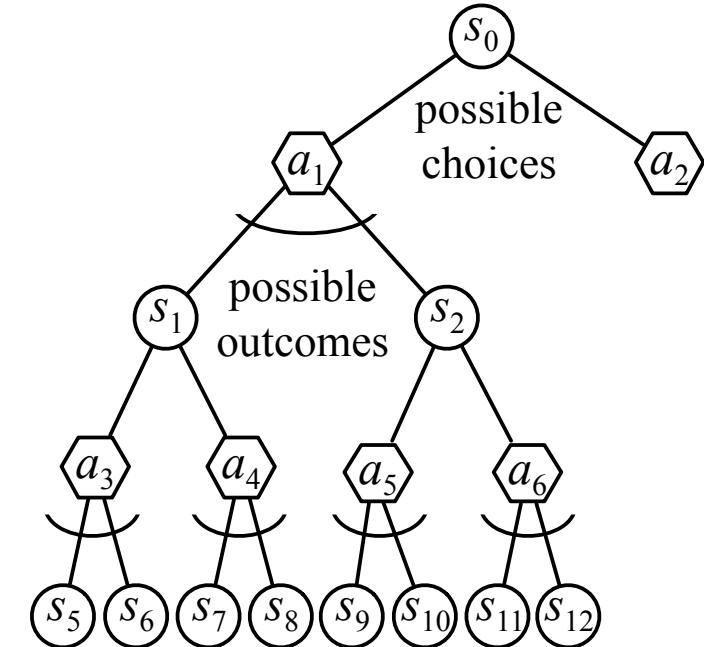
```

UCT( $s, h$ )
if  $s \in S_g$  then return 0
if  $h = 0$  then return  $V_0(s)$ 
if  $s \notin Envelope$  then do
    add  $s$  to  $Envelope$ 
     $n(s) \leftarrow 0$ 
    for all  $a \in Applicable(s)$  do
         $Q(s, a) \leftarrow 0; n(s, a) \leftarrow 0$ 
     $Untried \leftarrow \{a \in Applicable(s) \mid n(s, a) = 0\}$ 
    if  $Untried \neq \emptyset$  then  $\tilde{a} \leftarrow Choose(Untried)$ 
    else  $\tilde{a} \leftarrow argmin_{a \in Applicable(s)} \{Q(s, a) - C \times [\log(n(s)) / n(s, a)]^{1/2}\}$ 
     $s' \leftarrow Sample(\Sigma, s, \tilde{a})$ 
     $cost-rollout \leftarrow cost(s, \tilde{a}) + UCT(s', h - 1)$ 
     $Q(s, \tilde{a}) \leftarrow [n(s, \tilde{a}) \times Q(s, \tilde{a}) + cost-rollout] / (1 + n(s, \tilde{a}))$ 
     $n(s) \leftarrow n(s) + 1$ 
     $n(s, \tilde{a}) \leftarrow n(s, \tilde{a}) + 1$ 
return  $cost-rollout$ 

```

for $i - 1$ to n
 call $UCT(s_0, \infty)$
 $\forall s, \pi(s) \leftarrow argmin \{Q(s, a) \mid a \in Applicable(s)\}$

- As $n \rightarrow \infty$, π converges to optimal
 - ▶ Problem: finding optimal π may take many iterations



Using UCT Online

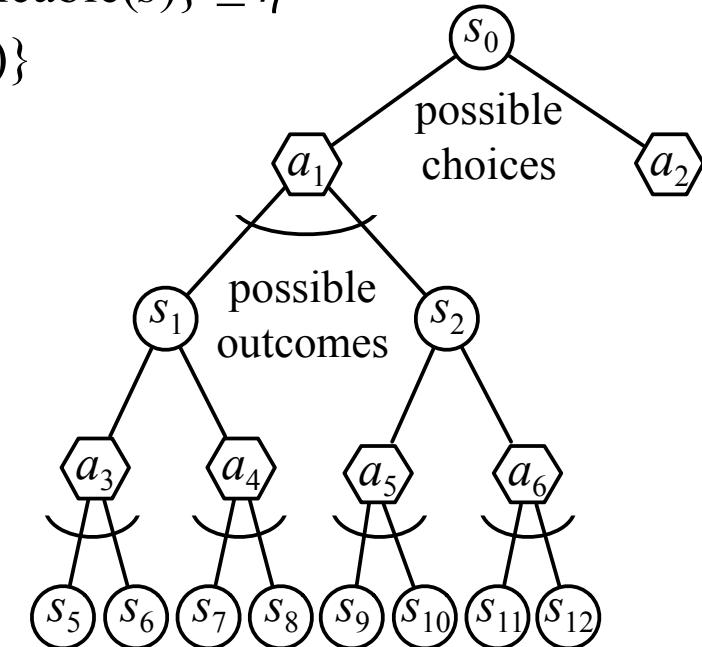
- Works better for online planning
 - ▶ $\pi(s_0)$ approaches optimal much faster than the rest of π
- Lookahead procedure for Run-Lookahead:
 - ▶ call $\text{UCT}(s, h)$ multiple times at current state s
 - ▶ e.g., until
 - allotted time runs out, or
 - $\max \{\text{last change in } Q(s,a) \mid a \in \text{Applicable}(s)\} \leq \eta$
 - ▶ return $\text{argmax } \{Q(s,a) \mid a \in \text{Applicable}(s)\}$

$\text{UCT}(s, h)$

```

if  $s \in S_g$  then return 0
if  $h = 0$  then return  $V_0(s)$ 
if  $s \notin \text{Envelope}$  then do
  add  $s$  to  $\text{Envelope}$ 
   $n(s) \leftarrow 0$ 
  for all  $a \in \text{Applicable}(s)$  do
     $Q(s, a) \leftarrow 0; n(s, a) \leftarrow 0$ 
   $Untried \leftarrow \{a \in \text{Applicable}(s) \mid n(s, a) = 0\}$ 
  if  $Untried \neq \emptyset$  then  $\tilde{a} \leftarrow \text{Choose}(Untried)$ 
  else  $\tilde{a} \leftarrow \operatorname{argmin}_{a \in \text{Applicable}(s)} \{Q(s, a) - C \times [\log(n(s)) / n(s, a)]^{\frac{1}{2}}\}$ 
   $s' \leftarrow \text{Sample}(\Sigma, s, \tilde{a})$ 
   $cost\text{-rollout} \leftarrow cost(s, \tilde{a}) + \text{UCT}(s', h - 1)$ 
   $Q(s, \tilde{a}) \leftarrow [n(s, \tilde{a}) \times Q(s, \tilde{a}) + cost\text{-rollout}] / (1 + n(s, \tilde{a}))$ 
   $n(s) \leftarrow n(s) + 1$ 
   $n(s, \tilde{a}) \leftarrow n(s, \tilde{a}) + 1$ 
return  $cost\text{-rollout}$ 

```



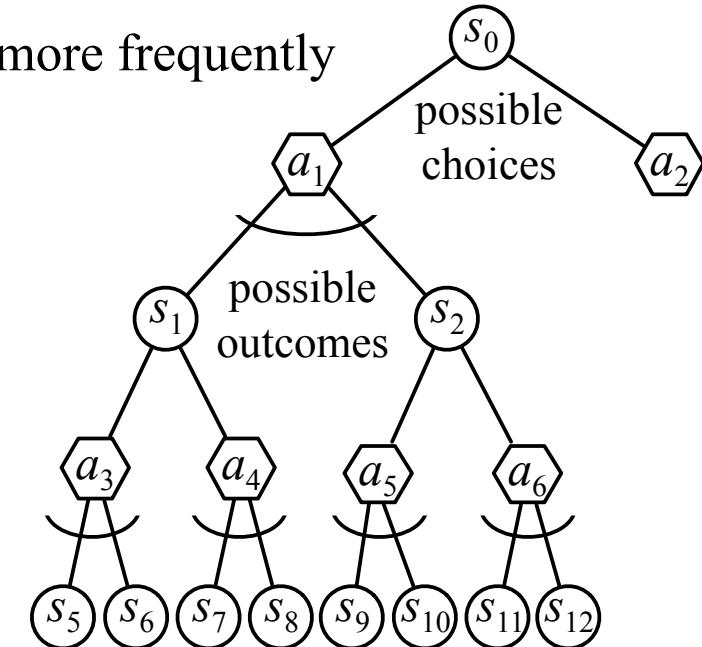
Using UCT Online

```

UCT( $s, h$ )
if  $s \in S_g$  then return 0
if  $h = 0$  then return  $V_0(s)$ 
if  $s \notin Envelope$  then do
    add  $s$  to  $Envelope$ 
 $n(s) \leftarrow 0$ 
for all  $a \in Applicable(s)$  do
     $Q(s, a) \leftarrow 0; n(s, a) \leftarrow 0$ 
 $Untried \leftarrow \{a \in Applicable(s) \mid n(s, a) = 0\}$ 
if  $Untried \neq \emptyset$  then  $\tilde{a} \leftarrow Choose(Untried)$ 
else  $\tilde{a} \leftarrow argmin_{a \in Applicable(s)} \{Q(s, a) - C \times [\log(n(s)) / n(s, a)]^{\frac{1}{2}}\}$ 
 $s' \leftarrow Sample(\Sigma, s, \tilde{a})$ 
 $cost-rollout \leftarrow cost(s, \tilde{a}) + UCT(s', h - 1)$ 
 $Q(s, \tilde{a}) \leftarrow [n(s, \tilde{a}) \times Q(s, \tilde{a}) + cost-rollout] / (1 + n(s, \tilde{a}))$ 
 $n(s) \leftarrow n(s) + 1$ 
 $n(s, \tilde{a}) \leftarrow n(s, \tilde{a}) + 1$ 
return  $cost-rollout$ 

```

- Lookahead procedure for Run-Lazy-Lookahead:
 - ▶ loop:
 - call UCT many times at current state
 - At state s visited, $\pi(s) \leftarrow$ action with highest expected utility
- Problem: the farther you follow π , the less likely that $\pi(s)$ is optimal
 - ▶ Near the bottom of the tree, $\pi(s)$ might be \approx random choice
- Possible workaround
 - ▶ Modify Run-Lazy-Lookahead to call UCT more frequently



Using UCT with a Simulator

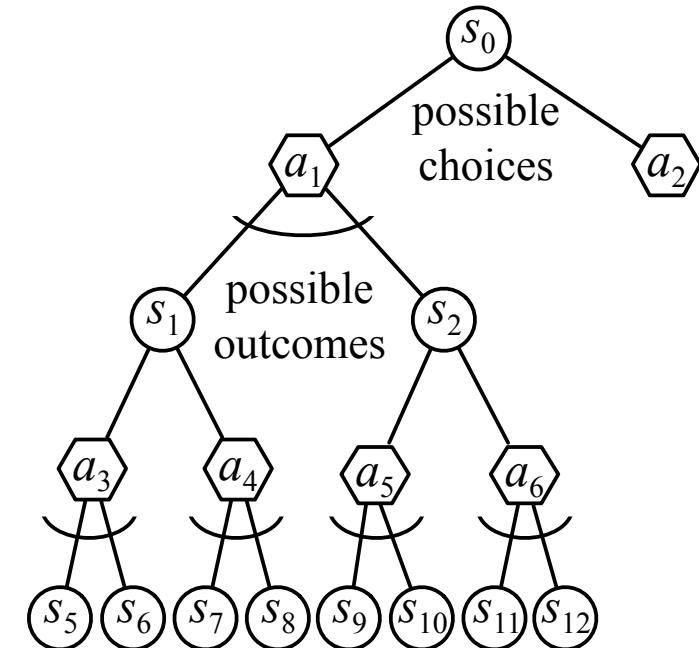
```

 $\text{UCT}(s, h)$ 
if  $s \in S_g$  then return 0
if  $h = 0$  then return  $V_0(s)$ 
if  $s \notin \text{Envelope}$  then do
    add  $s$  to  $\text{Envelope}$ 
     $n(s) \leftarrow 0$ 
    for all  $a \in \text{Applicable}(s)$  do
         $Q(s, a) \leftarrow 0; n(s, a) \leftarrow 0$ 
     $\text{Untried} \leftarrow \{a \in \text{Applicable}(s) \mid n(s, a) = 0\}$ 
    if  $\text{Untried} \neq \emptyset$  then  $\tilde{a} \leftarrow \text{Choose}(\text{Untried})$ 
    else  $\tilde{a} \leftarrow \operatorname{argmin}_{a \in \text{Applicable}(s)} \{Q(s, a) - C \times [\log(n(s)) / n(s, a)]^{\frac{1}{2}}\}$ 
     $s' \leftarrow \text{Sample}(\Sigma, s, \tilde{a})$ 
     $\text{cost-rollout} \leftarrow \text{cost}(s, \tilde{a}) + \text{UCT}(s', h - 1)$ 
     $Q(s, \tilde{a}) \leftarrow [n(s, \tilde{a}) \times Q(s, \tilde{a}) + \text{cost-rollout}] / (1 + n(s, \tilde{a}))$ 
     $n(s) \leftarrow n(s) + 1$ 
     $n(s, \tilde{a}) \leftarrow n(s, \tilde{a}) + 1$ 
return  $\text{cost-rollout}$ 

```

simulate a ; observe s'

- Suppose you don't know the probabilities and costs
 - ▶ But you have a fast, accurate simulator for the environment
- Lookahead procedure (see previous slides)
 - ▶ Run UCT multiple times in the simulated environment
 - ▶ Learn state-transition probabilities, expected utilities



Using UCT for Exploration

$\text{UCT}(s, h)$

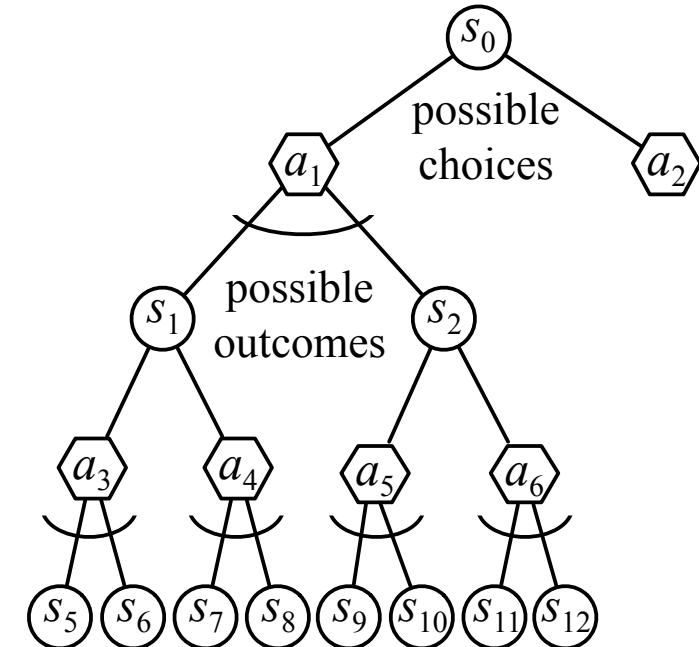
```

if  $s \in S_g$  then return 0
if  $h = 0$  then return  $V_0(s)$ 
if  $s \notin \text{Envelope}$  then do
    add  $s$  to  $\text{Envelope}$ 
     $n(s) \leftarrow 0$ 
    for all  $a \in \text{Applicable}(s)$  do
         $Q(s, a) \leftarrow 0$ ;  $n(s, a) \leftarrow 0$ 
     $\text{Untried} \leftarrow \{a \in \text{Applicable}(s) \mid n(s, a) = 0\}$ 
    if  $\text{Untried} \neq \emptyset$  then  $\tilde{a} \leftarrow \text{Choose}(\text{Untried})$ 
    else  $\tilde{a} \leftarrow \operatorname{argmin}_{a \in \text{Applicable}(s)} \{Q(s, a) - C \times [\log(n(s)) / n(s, a)]^{\frac{1}{2}}\}$ 
     $s' \leftarrow \text{Sample}(\Sigma, s, \tilde{a})$ 
     $\text{cost-rollout} \leftarrow \text{cost}(s, \tilde{a}) + \text{UCT}(s', h - 1)$ 
     $Q(s, \tilde{a}) \leftarrow [n(s, \tilde{a}) \times Q(s, \tilde{a}) + \text{cost-rollout}] / (1 + n(s, \tilde{a}))$ 
     $n(s) \leftarrow n(s) + 1$ 
     $n(s, \tilde{a}) \leftarrow n(s, \tilde{a}) + 1$ 
return  $\text{cost-rollout}$ 

```

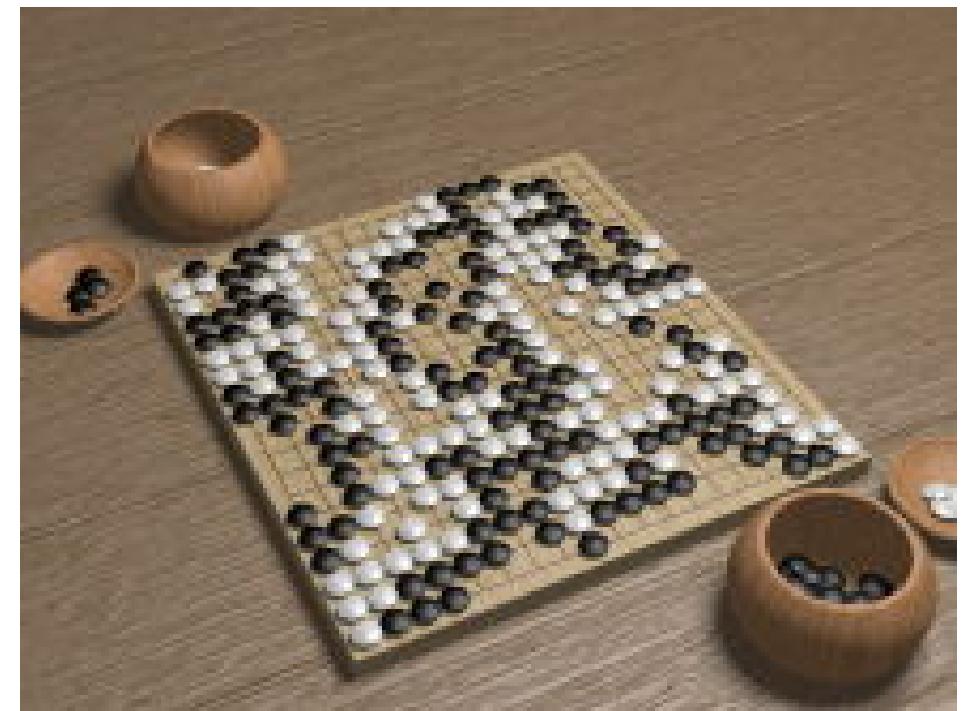
perform a ; observe s'

- Suppose you don't know the probabilities and costs
- Suppose you can restart your actor as many times as you want
 - ▶ Caveat: usually not very feasible in real environments
- Can modify UCT to be an acting procedure, use it to explore the environment



UCT in Two-Player Games

- Generate Monte Carlo rollouts using a modified version of UCT
- Main differences:
 - ▶ Instead of choosing actions that minimize accumulated cost, choose actions that maximize payoff at the end of the game
 - ▶ UCT for player 1 recursively calls UCT for player 2
 - Choose opponent's action
 - ▶ UCT for player 2 recursively calls UCT for player 1
- First competent computer programs for go
 - ▶ ≈ 2008–2012
- Monte Carlo rollout techniques similar to UCT were used to train AlphaGo



Summary

- SSPs
- solutions, closed solutions, histories
- unsafe solutions, acyclic safe solutions, cyclic safe solutions
- expected cost, planning as optimization
- policy iteration
- value iteration (asynchronous version)
 - ▶ Bellman-update
- AO*, LAO*
- Planning and Acting
 - ▶ Run-Lookahead
 - ▶ FS-Replan
- UCB, UCT