Markov? memoriless
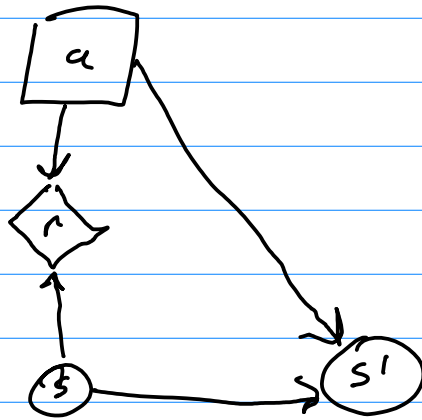$$P(x_{t+1} | x_0 .. x_t) = P(x_{t+1} | x_t)$$

Today
What is an MDP?
What is a ~~policy~~?
How to evaluate policies
Policy Search

## MDPs



$$\text{maximize} \sum_{t=0}^{\infty} \gamma^t r_t$$

## "Tuple Definition"

$$(S, A, R, T, \gamma) + p_0$$

$s \in S$

$s = (x, y)$  $S = \mathbb{R}^4$

$S$ - state space - set of states $\left( \begin{array}{c} \text{e.g. } \{1, 2, 3\}, \ \mathbb{R}^4 \\ \{\text{working, malfunctioning}\} \\ [0, 1]^4 \end{array} \right.$

$A$ - action space - set of actions

$$A(s)$$

$R$ - reward function   $R : S \times A \times S \to \mathbb{R}$

$$R(s, a) \equiv \underset{s'}{E}\left[ R(s, a, s') \right]$$

$T$ - "transition kernel"

   Explicit   or   Implicit   ("Generative Model")

$$T(s' | s, a) \qquad s' = G(s, a)$$

$\gamma$ - discount     $\gamma \in [0, 1)$

$p_0$ - initial state distribution

---

## Breakout Rooms

Cooking a pot of Pasta

$$(S, A, R, T, \gamma)$$

Team 5

$$S = \{1 \ldots 10\} \times \{1 \ldots 10\} \times \{0 \ldots 5\}$$

water temp       softness       saltiness

$A = \{up, down, add\ salt, eat, taste, add\ sauce,\}$

$R = + tastiness^{-(softness-5)^2} - 1$ time step
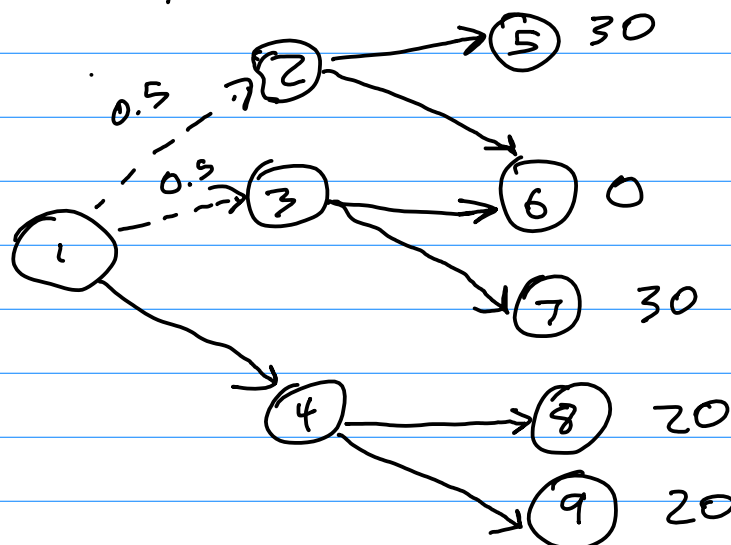
$T =$ if $a = high$ temp increases, if temp $= 10$
                                        softness $+= 1$

$\gamma = 0.99$

## Policies   determine what actions are taken

Open Loop:   sequence
Closed Loop:   $\pi: S \rightarrow A$



$A = \{up, down\}$

open loop policies

$(\uparrow, \uparrow)$  $0.5 \cdot 30 = 15$
$(\uparrow, \downarrow)$    $15$
$(\downarrow, \downarrow)$    $20$
$(\downarrow, \uparrow)$    $20$

Closed Loop $\left(\uparrow, \begin{cases} \uparrow & if\ 2 \\ \downarrow & ow \end{cases}\right)$   $0.5 \cdot 30 +$
                                                                                                      $0.5 \cdot 30 = 30$

## Evaluation

estimate $\overset{return}{u} = \underset{s_0 \sim p_0}{E}\left[\sum_{t=0}^{\infty} \gamma^t r_t \mid a_t = \pi(s_t)\right]$

$$\hat{u} = \sum_{t=0}^{T-1} \gamma^t r_t$$

**Simulation**

$s \leftarrow \text{sample}(\beta_0)$

$\hat{u} \leftarrow 0$

for $t$ in $0 \ldots T-1$

$\qquad s', r \leftarrow G(s,a) \quad\leftarrow a \leftarrow \pi(s)$

$\qquad \hat{u} += \gamma^t r$

$\qquad s \leftarrow s'$

return $\hat{u}$

$$u \approx \bar{u}_m = \frac{1}{m} \sum_{i=1}^{m} \hat{u}_i \quad\leftarrow \text{from simulation}$$

$\qquad\qquad \uparrow$ Monte Carlo Estimate