

MDP

S : {states}

A : {Actions}

T : transition function $T(s_t, a, s_{t+1}) = \Pr(s_{t+1} | s_t, a)$ PDF over states at time $t+1$

R : reward function

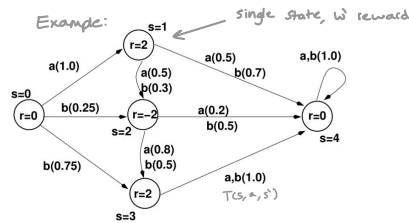
↳ Simple case: $R(s)$ fixed for a given state
 ↳ complex case: $R(s, a, s')$

At each time step t , agent is in some state, s_t and must take an action, a_t .
Each action causes a transition to a new state, s_{t+1} .

Representations:

S-A-R Space:

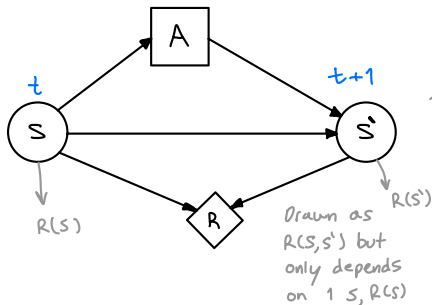
more compact



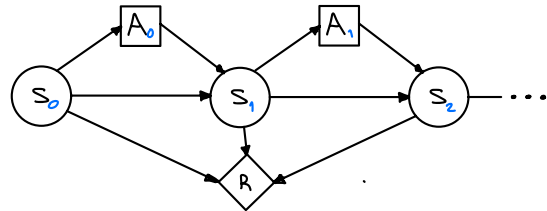
DBN:

* technically a DBN, since R & A aren't R.V's

Compact: represent each state w only 2 time slices



Unrolled (Full DBN)



For as many time steps as we like

Optimal policy π^* , gives max expected reward. For $t \rightarrow \infty$: $\sum_{t=0}^{\infty} \gamma^t R(s_t)$

"Value" of being in s w t stages to go, $V(s, t)$.

↳ Find w DP: for all s' what action gets us to best-next state?

Start w $V^0(s) = R(s)$

in practice, until V^t stops changing much

Then, w t -stages to go

∞ stage to go

Optimal

$$V^t(s) = \max_a [R(s) + \gamma \sum_{s'} \Pr(s' | s, a) V^{t-1}(s')]]$$

$$V^*(s)$$

$$\pi^t(s) = \arg \max_a V^t(s)$$

$$\pi^*(s)$$