

Model 1: predict what env will do next

$$\text{Value func: } V_\pi(s) = E_\pi[R_{t+1} + \gamma R_{t+2} + \dots | S_t = s]$$

prediction of future reward

$$\text{Reward (R): } R_s = E[R_{t+1} | S_t = s]$$

$$\text{Return (G): } G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Markov reward process  $\langle S, P, R, \gamma \rangle$

Bellman eq: immediate reward + successor reward

$$V(s) = E[R_{t+1} + \gamma V(S_{t+1}) | S_t]$$

$$= R + \gamma P V_{t+1} = (I - \gamma P)^{-1} R$$

$O(\# \text{ of states}^3)$

DP

MC

TD

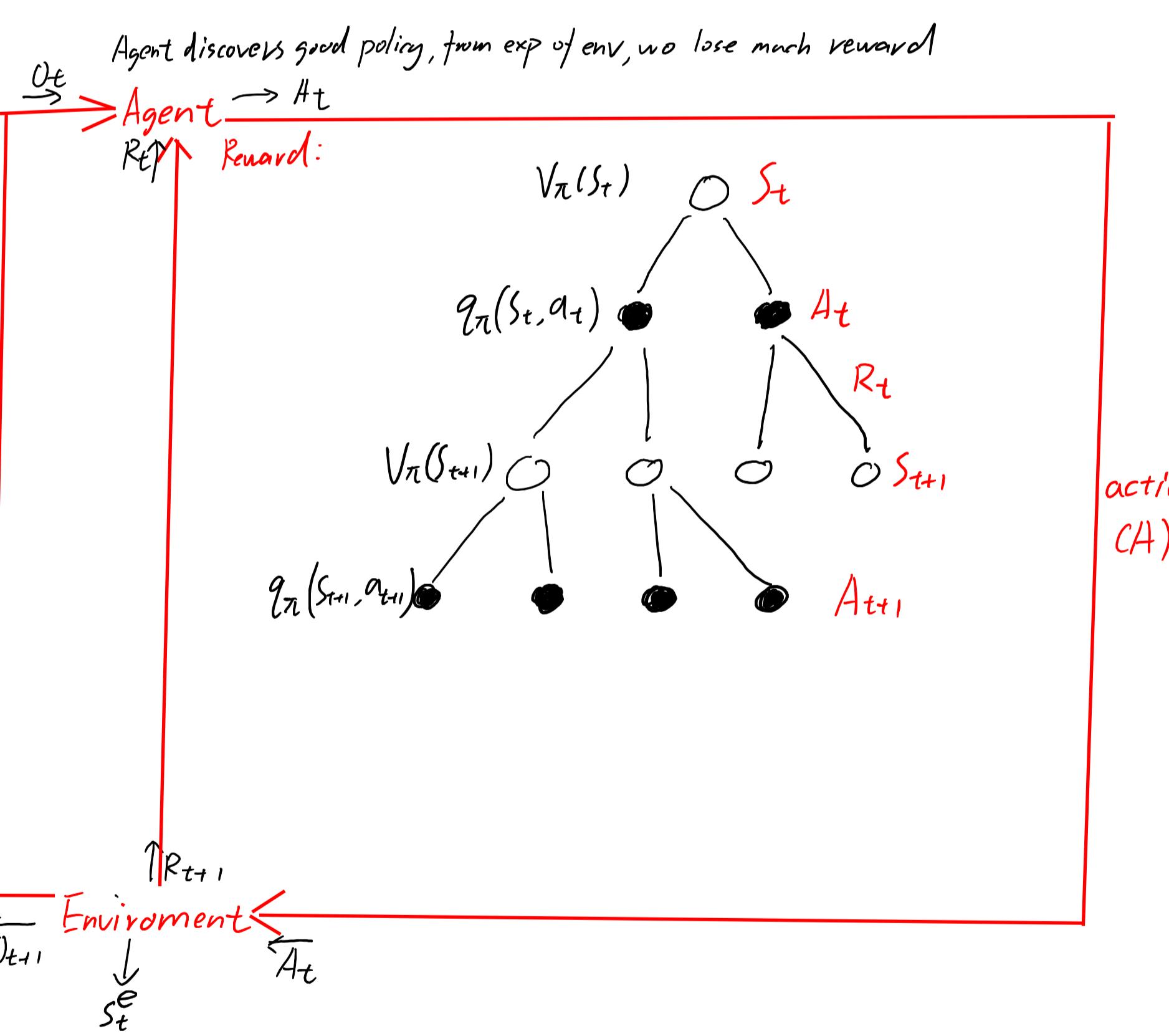
infinite states  $\Rightarrow$  policy gradient

$$\text{State trans prob (P): } P(s'|s, a)$$

state trans matrix:  $[\# \text{ of states} \times \# \text{ of states}]$

$$\text{Markov state (S)}$$

$$P(s_{t+1}|s) = P(s_{t+1}|s, \dots, s_t)$$



$$\begin{aligned} \text{Full: } O_t &= S_t^o - S_t^e \Rightarrow MDP \\ \text{Partial: } S_t^o &\neq S_t^e \xrightarrow{\text{conversion}} \end{aligned}$$

$$\text{Policy } (\pi): \text{map from state to action; deterministic: } a = \pi(s) \text{ stochastic: } \pi(a|s) = P(A_a|S_s) \text{ strict sense stationary (time independent)}$$

Markov decision process  $\langle S, P, R, \gamma, A \rangle$

Only depends on current state

$$P_{ss'}^\pi = \sum_{a \in A} \pi(a|s) P_{ss'}^a, \quad R_s^\pi = \sum_{a \in A} \pi(a|s) R_s^a$$

State value func:

$$\begin{aligned} \text{expected return following policy } \pi \\ V_\pi(s) = E_\pi[G_t | S_t] &= E_\pi\left[R_{t+1} + \gamma V_\pi(S_{t+1}) | S_t\right] \\ &= \sum_{a \in A} \pi(a|s) \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_\pi(s') \right) \end{aligned}$$

Action value func:

$$\begin{aligned} \text{expected return following action then policy} \\ q_\pi(s, a) &= E_\pi[G_t | S_t, A_t] = E_\pi\left[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) | S_t, A_t\right] \\ &= R_s^a + \gamma \sum_{s' \in S} \sum_{a' \in A} \pi(a'|s') q_\pi(s', a') \end{aligned}$$