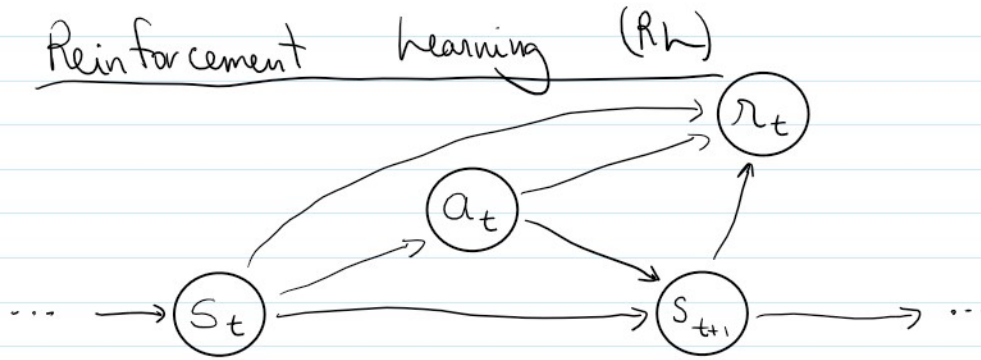


Tutorial #12

December 2, 2021 1:00 PM



- An agent:
 - 1) begins in state s_t
 - 2) takes an action a_t
 $a_t \sim \text{policy } \pi$
 - 3) moves to the next state s_{t+1}
transition probs. $P[s_{t+1} | s_t, a_t]$
 - 4) receives a reward/cost r_t

• An episode of length T :

$$\tau := (s_0, a_0, r_0, s_1, \dots, a_{T-1}, r_{T-1}, s_T).$$

- Goal of RL: find an optimal **policy** π that maximizes the value (or Q) function.

$$V^a(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right].$$

• Suppose we have a randomized policy π^θ .

$$P(\tau \mid \text{randomized } \pi^\theta) \\ = \underbrace{P(s_0)}_{\text{initial state}} \cdot \prod_{t=0}^{T-1} \underbrace{\pi^\theta(a_t | s_t)}_{\text{policy for } a_t} \cdot \underbrace{P[s_{t+1} | s_t, a_t]}_{\text{prob. for } s_{t+1}}$$

Then $\nabla_\theta \log P(\tau \mid \pi^\theta)$

$$= \nabla_\theta \left[\log P(s_0) + \sum_{t=0}^{T-1} \log \pi^\theta(a_t | s_t) + \log (P[s_{t+1} | s_t, a_t]) \right] \\ = \sum_{t=0}^{T-1} \nabla_\theta \log \pi^\theta(a_t | s_t).$$

• Likelihood-ratio trick:

$$P[x; \theta] \nabla_\theta \log P[x; \theta] = \nabla_\theta P[x; \theta].$$

• Bellman Principle:

$$V(s) = \max_{a \in \mathcal{A}} \mathbb{E} \left[r(a, s, s') + \gamma V(s') \right]$$

• Q-function:

$$Q(s, a) = \mathbb{E} \left[r(a, s, s') + \gamma V(s') \right] \\ = \mathbb{E} \left[r(a, s, s') + \gamma \max_{a' \in \mathcal{A}} Q(s', a') \right]$$