

Bayesian Reinforcement Learning

RL: $\langle S, A, R \rangle$

we are missing T (transition dynamics)

MDP: $\langle S, A, T, R \rangle \Rightarrow \Rightarrow \pi^*$
 $\Rightarrow V^*$

Bayesian Q-learning:

distributions over V^*

Model based Bayesian Exploration

distributions over T

Model-based Bayesian RL:

- distributions over $P_n(s'|s, a)$

- Let $\Theta_{s, s'}^a = P_n(s'|s, a)$

Prior $P(\Theta_{s, s'}^a) = P_n(P_n(s'|s, a) = V)$

Posterior $P(\Theta_{s, s'}^a | s \xrightarrow{a} s')$

Dirichlet distributions are conjugate priors of multinomial distributions

$$P_n(x_1, \dots, x_n) = k \Theta_1^{\alpha_1 - 1} \Theta_2^{\alpha_2 - 1} \dots \Theta_m^{\alpha_m - 1}$$

$$\text{Let } P(\theta_{s,s'}^a) = \text{Dir}(\alpha_{s,s'}^a)$$

$$\text{then } P(\theta_{s,s'}^a | s \xrightarrow{a} s') = R P(s \xrightarrow{a} s' | \theta_{s,s'}^a) P(\theta_{s,s'}^a)$$

$$= R P(s \xrightarrow{a} s' | \theta_{s,s'}^a) \prod_{s''} (\theta_{s,s'}^a)^{\alpha_{s,s'}^a - 1}$$

$$= R \theta_{s,s'}^a \prod_{s''} (\theta_{s,s''}^a)^{\alpha_{s,s''}^a - 1}$$

$$= R \prod_{s'' \neq s'} (\theta_{s,s''}^a)^{\alpha_{s,s''}^a - 1} (\theta_{s,s'}^a)^{\alpha_{s,s'}^a}$$

$$= \text{Dir}(\alpha_{s,s'}^a + \delta(s', s''))$$

State space: $S \times \{ \theta_{s,s'}^a \}$

↓
Bellman
steps
numbers

Hence we have a POMDP