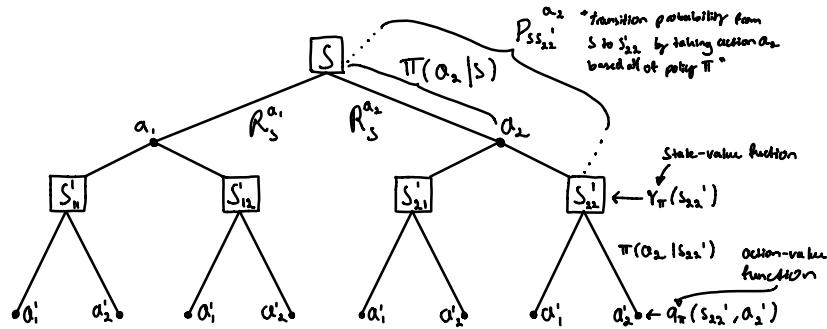


Inference

Markov Decision Process: $\langle S, A, P_{ss'}, R_s^a, \pi^a, \gamma \rangle$
 graph



$$V_\pi(s) = \sum_{a \in A} \pi(a|s) [R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_\pi(s')]$$

$$q_\pi(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \cdot \sum_{a' \in A} \pi(a'|s') \cdot q_\pi(s', a')$$

Ex:

$$V_\pi(s) = \pi(a_1|s) \cdot [R_s^{a_1} + \gamma \left[P_{SS'_{11}}^{a_1} \cdot Y_\pi(S'_{11}) + P_{SS'_{12}}^{a_1} \cdot Y_\pi(S'_{12}) \right]] + \pi(a_2|s) \cdot [R_s^{a_2} + \gamma \left[P_{SS'_{21}}^{a_2} \cdot Y_\pi(S'_{21}) + P_{SS'_{22}}^{a_2} \cdot Y_\pi(S'_{22}) \right]]$$

probability we go from \$s \rightarrow s'\$ scaled by the value at the state scaled by the likelihood we end up there

multiplied by the value of \$S'_{11}\$

discounted sum of value of branches scaled by the likelihood we end up at one at the branches

Value of immediate transition by \$a\$, plus averaged value of the states we would end up in given that we took \$a\$ at \$s\$

Sum of probability we end up at a state multiplied by its value

> the whole thing:

The value of a state is defined as the average value of taking

Starting boundary case:

[going from $V \rightarrow V_\pi$]

1b.

$$q_\pi(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \cdot V_\pi(s')$$

the value at taking an action, a , at a state, s , under policy, π , represented by $q_\pi(\cdot, \cdot)$ is the reward for taking the action, R_s^a , which is given by the environment, in addition to the discounted [γ] average value of the states you would end up at.

determined by env.
odds we'll end up there - value
- high value - high odds → big influence on q_π
- low value - high odds →

1a.

$$V_\pi(s) = \sum_{a \in A} \pi(a|s) \cdot q_\pi(s, a)$$

$\pi(a|s)$

assume a stationary policy, one that doesn't change throughout an episode

Value, $V_\pi(\cdot)$ of state, s , under policy, π , is the average value of the actions, $q_\pi(s, \cdot)$ available

the intuition is that an action leads you to another state, which itself respects the future average value

as a is more valuable it'll lead you to more valuable states

- By averaging, we take into account probability we will pick an action/odd well end up somewhere as well as true value of the result →
- Both of these definitions are mindful of the fact that taking an action does not translate to going to another state.
- This notion also explains why reward, R_s^a , comes from the result of taking a , not the intent of taking a .

1a. & 1b. define quantities by looking a single step ahead (1 lego brick graph). By looking further we can form recursive definitions.

2b.

$$q_\pi(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \left[\sum_{a' \in A} \pi(a'|s) q_\pi(s', a') \right]$$

- the value of an action available at state s is the sum of the immediate reward, and the discounted average of the states we'll end up at, which is represented as the average value of the actions available at those future states

we value an action by its immediate reward and the various actions we'll take as a consequence at the current one (which includes the rewards of future actions)

1a./b. lead us to 2a./b. which describe methods for navigating an MDP under policy π . It's practical to figure out an optimal policy, π^* for solving problems posed as MDPs.

best possible performance in an MDP

$$V^*(s) = \max_{\pi} V_{\pi}(s) \leftarrow \text{the optimal state-value function, } V^*(\cdot), \text{ is the one where it has highest value across all possible policies}$$

$$Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a) \leftarrow \text{the optimal action-value function, } Q^*(\cdot, \cdot), \text{ is the best action-value across all policies}$$

scalar comparison

with this definition, can optimize π w/ $Q \gg V$ as objective functions

assume $\pi_1 \geq \pi_2 \Rightarrow V_{\pi_1}(s) \geq V_{\pi_2}(s)$

- there exists optimal policy, $\pi^* \geq \pi, V^* \geq V_{\pi}$
- for clarity, $V_{\pi^*}(s) = V^*(s) \quad \& \quad Q_{\pi^*}(s, a) = Q^*(s, a)$

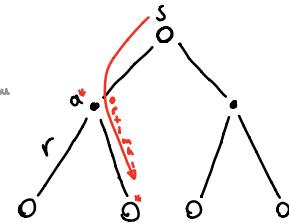
since you've always made the right choice

3b.

$$q_{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a q_{\pi}(s')$$

continue to make best local selections

$$= R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \left[\max_a q_{\pi}(s', a) \right]$$



3a.

$$\begin{aligned} V_{\pi}(s) &= \max_a q_{\pi}(s, a) \\ &= \max_a \left[R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi}(s') \right] \end{aligned}$$

The addition at "openn" (3a/3b.) only modifies the binary/localized equations, and does not keep them in their optimal form.

To this aligns w/ the idea that for an optimal outcome we need to make a locally good choice, and obtain well more a good choice in the future too.

Bellman equality is non-linear

- ↳ these are the basis for
 - Value iteration
 - Policy iteration
 - Q-learning
 - Sarsa