

MRP is a tuple $(S, A, P, R, \gamma, \pi)$, where S is the state space A is the action space $P_{ss'}^a$ is the transition probability $p(s_{t+1}|a, s_t)$ $R_{ss'}^a$ is a reward function defined as $r_s = r(s, a, s') = r(s_t, s_{t+1}, a)$, that is collected upon leaving a the state s , i.e. at time step $t + 1$ γ is the discount factor π policy which determines what action should be taken $\pi(a|s)$ - Stochastic: the action is distributed over the policy $a \sim \pi(a|s)$ - Deterministic: the action is set to the result of the policy $a = \pi(a|s)$

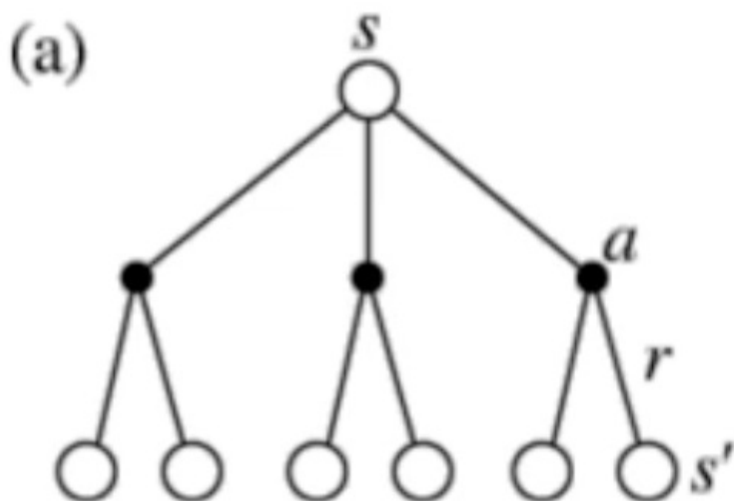
Value of a state

The value of a state V^π given a certain policy π will be followed thereafter, is the expected discounted return after following the policy: $V^\pi(s) = \mathbb{E}_{\pi} [R_t | S_t = s]$

$$\begin{aligned} &= \mathbb{E}_{\pi} \left[r_{t+1} + \sum_{k=1}^{\infty} \gamma^k r_{t+k+1} \middle| S_t = s \right] \\ &= \mathbb{E}_{\pi} \left[r_{t+1} + \gamma \sum_{k=1}^{\infty} \gamma^{k-1} r_{t+k+1} \middle| S_t = s \right] \\ &= \mathbb{E}_{\pi} \left[r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \middle| S_t = s \right] \\ &= \mathbb{E}_{\pi} [r_{t+1} | S_t = s] + \gamma \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \middle| S_t = s \right] \end{aligned} \quad \text{Equation (1)}$$

The first equation considers the expected immediate award given we are in state s . In other words the question is. Note the difference between the return R_t and the immediate reward r_t

The immediate reward can be easily calculated using a backup tree like this one:



- The root of the tree denotes the current state s .
- The tree branches off into every possible state s' that we can end up at from state s The question now is what is the probability to take a certain action?
- This is given by the policy $\pi(a|s)$

Next question is what is the probability that we end up in a state s' given an action a is taken? (The model might perform an action, but the action might not succeed 100% of the time, e.g. try to move left, but the environment pushed the agent forward instead) - This is given by the transition probability matrix $P_{ss'}^a = P(s, a, s')$

Therefore equation (1) could be expressed as:

$$\mathbb{E}_{\pi} [r_{t+1} | S_t = s] = \sum_{a \in A} \pi(a|s) \sum_{s' \in S} P(s'|s, a) r(s, a, s')$$

The expression uses the immediate when moving to the next state. This part of the equation is concerned only with what reward does the agent get rather than the value of the state

Equation (2)

$$\gamma \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \middle| S_t = s \right]$$

- The γ has been moved outside as the expected value does not change by constants - This equation looks very much like the value of the state s' , but the difference is that it is conditioned on $S_t = s$ **and not** on $S_t = s'$ like shown below:

$$V^\pi(s') = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \middle| S_t = s' \right]$$

- In order to calculate (2) a similar approach needs to be taken as in equation (1) with the backup tree incorporating the dynamics of the environment: the policy $\pi(a|s)$ and the transition probability matrix $P(s'|a, s)$

$$\begin{aligned} \gamma \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} | S_t = s \right] &= \sum_{a \in A} \pi(a|s) \sum_{s' \in S} P(s'|s, a) \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \\ &= \sum_{a \in A} \pi(a|s) \sum_{s' \in S} P(s'|s, a) \gamma V^{\pi}(s') \end{aligned}$$

- After incorporating those, the last term $\sum_{k=0}^{\infty} \gamma^k r_{t+k+2}$, denotes the expected discounted return from state $t + 1$, i.e. $V^{\pi}(s')$ ##### Combining the two \$\$

$$\begin{aligned} V^{\pi}(s) &= \mathbb{E}_{\pi}[R_t | S_t = s] \\ &= \mathbb{E}_{\pi}[r_{t+1} + \gamma \sum_{k=1}^{\infty} \gamma^{k-1} r_{t+k+1} | S_t = s] \\ &= \sum_{a \in A} \pi(a|s) \sum_{s' \in S} P(s'|s, a) (r(s, a, s') + \gamma \sum_{k=1}^{\infty} \gamma^{k-1} r_{t+k+1}) \\ &= \sum_{a \in A} \pi(a|s) \sum_{s' \in S} P_{ss'}^a (R_{ss'}^a + \gamma V^{\pi}(s')) \end{aligned}$$

\$\$

Policy Evaluation

- The value function $V(s)$ creates an ordering of all the policies π .
- An optimal policy always exists (sometimes > 1)
- An optimal policy needs to yield returns that are greater or equal to any other returns generated by a different policy. In other words - the policy would find the best solution no matter from the starting state

$$V^*(s) = \max_{\pi} V^{\pi}(s), \forall s \in S$$