

- Markov decision processes formally describe an environment for reinforcement learning.
- Where the environment is fully observable.
- i.e. The current state completely characterizes the process.
- Almost all RL problems can be formalized as MDPs e.g.
 - Optimal control primarily deals with continuous MDPs
 - Partially observable problems can be converted to MDPs
 - Bandits are MDPs with one state.

Markov Property

The future is independent of the past given the present.

Def: The state S_t is Markov iff -

$$P[S_{t+1} | S_t] = P[S_{t+1} | S_1, \dots, S_t]$$

The state captures all relevant information from history. Once the state is known, history maybe thrown away i.e., state is a sufficient statistic of the future.

State Transition Matrix:

For a Markov state s and a successor state s' , the state transition probability is defined by -

$$P_{ss'} = P[S_{t+1} = s' | S_t = s]$$

State transition matrix P defines transition probabilities from all states s to all successor states s' .

$$P = \text{from } \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix} \text{ to } \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix}$$

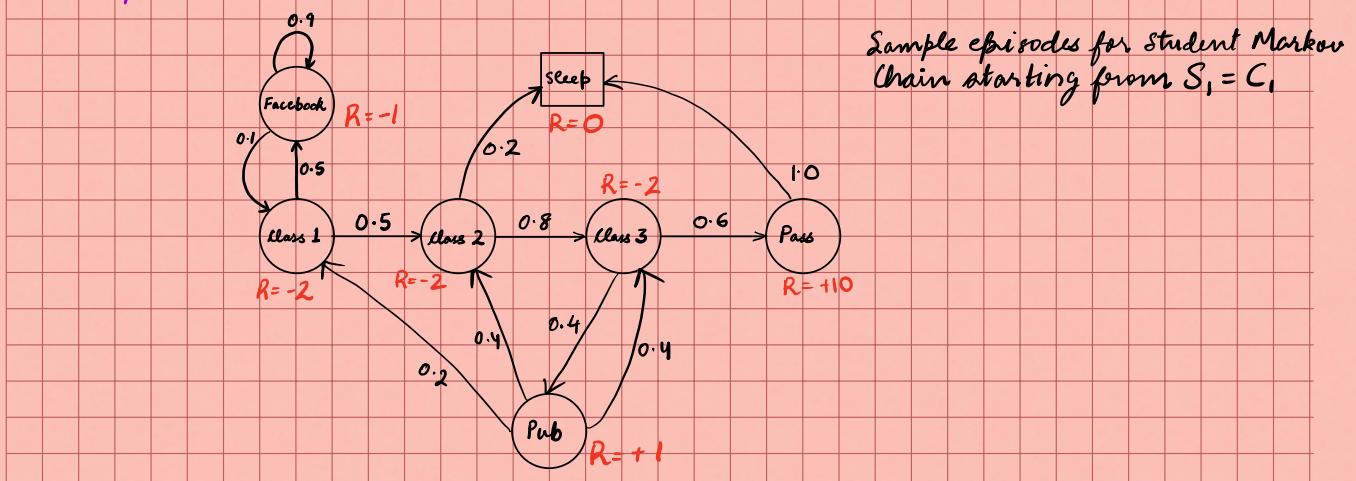
where each row of the matrix sums to 1.

Markov Process - It is a memoryless random process, i.e., a sequence of random states S_1, S_2, \dots with Markov property.

Def: - A Markov process or a Markov chain is a tuple $\langle S, P \rangle$:

- S is a finite set of states
- P is a state transition probability matrix, $P_{ss'} = P[S_{t+1} = s' | S_t = s]$

Example: Student Markov Chain:



Markov Reward Process: It is a Markov chain with values.

Def: Markov reward process is a tuple $\langle S, P, R, \gamma \rangle$

- S is a finite set of states

- P is a state transition probability matrix: $P_{ss'} = P[S_{t+1} = s' | S_t = s]$

- R is a reward function, $R_s = E[R_{t+1} | S_t = s]$

- γ is a discount factor, $\gamma \in [0, 1]$

Def: The return G_t is the total discounted reward from time-step t .

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

The discount $\gamma \in [0, 1]$ is the present value of future rewards.

Why discount?

- There's more uncertainty in the future.

- Mathematical convenience - avoids infinite rewards in cyclic processes

- Sometimes undiscounted rewards are used if all sequences terminate.

Value Function

Def: The state value function $v(s)$ of an MRP is the expected return starting from state s .

$$v(s) = E[G_t | S_t = s]$$

Bellman Equation for MRPs

The value function can be decomposed into two parts.

- immediate reward R_{t+1}

- discounted value of successor state $\gamma v(S_{t+1})$

$$\begin{aligned} v(s) &= E[G_t | S_t = s] \\ &= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s] \\ &= E[R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots | S_t = s)] \\ &= E[R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= E[R_{t+1}] + \gamma E[G_{t+1} | S_t = s] \\ &= R_s + \gamma \sum_{s' \in S} P_{ss'} v(s') \end{aligned}$$

Bellman eq. can be written in matrix form.

$$\begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix} = \begin{bmatrix} R_1 \\ \vdots \\ R_n \end{bmatrix} + \gamma \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix} \begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix}$$

$$v = R + \gamma P v$$

$$(I - \gamma P) v = R$$

$$v = (I - \gamma P)^{-1} R$$

- Computational complexity $\sim O(n^3)$

- Direct solution only feasible for small MRPs
- More iterative processes to solve them
 - o Dynamic Programming
 - o Monte-Carlo Simulation
 - o Temporal-Difference learning

Markov Decision Process:

A Markov Decision Process (MDP) is a MRP with decisions. It is an environment in which all states are Markov.

Definition: A MDP is a tuple $\langle S, A, P, R, \gamma \rangle$ where :

- S is a finite set of states
- A is a finite set of actions.
- P is a state transition probability matrix.

$$- P_{ss'}^a = P[S_{t+1} = s' | S_t = s, A_t = a]$$

$$- R_s^a = E[R_{t+1} | S_t = s, A_t = a]$$

$$- \gamma \text{ is a discount factor } \gamma \in [0, 1]$$

Policy: A policy π is a distribution over actions given states,

$$\pi(a | s) = P[A_t = a | S_t = s]$$

A policy fully defines the behaviour of the model

- MDP policies depend on the current state and not the history.

- Policies are stationary (time-independent)

$$A_t \sim \pi(\cdot | S_t), \forall t > 0$$

Given an MDP $M = \langle S, A, P, R, \gamma \rangle$ and a policy π .

The state sequences S_1, S_2, S_3, \dots is a Markov Process $\langle S, P^\pi \rangle$

The state and reward sequence S_1, R_2, S_2, \dots is a MRP $\langle S, P^\pi, R^\pi, \gamma \rangle$

where -

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

$$R_s^\pi = \sum_{a \in A} \pi(a | s) R_s^a$$

Value Function

The state-value function $v_\pi(s)$ of an MDP is the expected return starting from state s , and then following policy π .

$$v_\pi(s) = E_\pi [G_t | S_t = s]$$

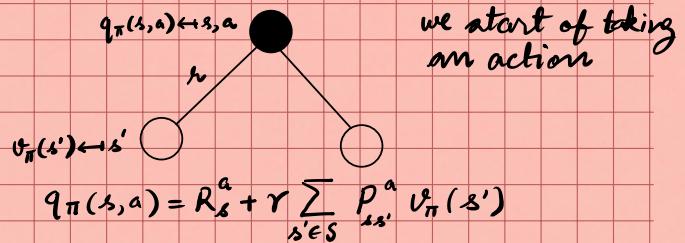
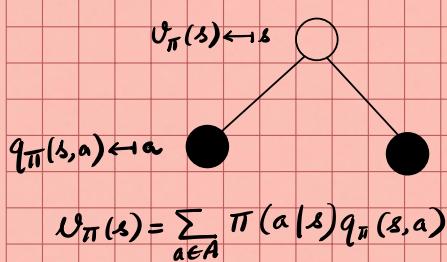
The action-value function $q_{\pi}(s, a)$ is the expected returns starting from state s and taking action a , and then following π .

$$q_{\pi}(s, a) = E_{\pi}[G_t | S_t = s, A_t = a]$$

Bellman Expectation Equation:

$$V_{\pi}(s) = E_{\pi}[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s]$$

$$q_{\pi}(s, a) = E_{\pi}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$$



avg. over actions that we might take

V : how good is it to be in a particular state

q : how good is it to take a particular action

From both equations:

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

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

Optimal Value Function:

Def: The optimal state value function is the max. value function over all the policies.

$$V_*(s) = \max_{\pi} V_{\pi}(s)$$

The optimal action-value function $q_*(s, a)$ is the maximum action-value function over all policies.

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

The optimal value function specifies the best possible performance in the MDP.

An MDP is "solved" when we know the optimal value function.

Optimal Policy:

Define a partial ordering over policies:

$$\Pi \geq \Pi' \text{ if } V_{\Pi}(s) \geq V_{\Pi'}(s), \forall s$$

Theorem:

For any Markov Decision Process -

- There exists an optimal policy Π^* that is better than or equal to all other policies, $\Pi^* \geq \Pi, \forall \Pi$
- All optimal policies reach the optimal value function $V_{\Pi^*}(s) = V^*(s)$
- All optimal policies achieve the optimal action-value function $q_{\Pi^*}(s, a) = q^*(s, a)$

Finding an Optimal Policy -

An optimal policy can be found by maximizing over $q^*(s, a)$,

$$\Pi^*(a|s) = \begin{cases} 1 & \text{if } a = \underset{a \in A}{\operatorname{argmax}} q^*(s, a) \\ 0 & \text{otherwise} \end{cases}$$

- There is always a deterministic optimal policy for any MDP.