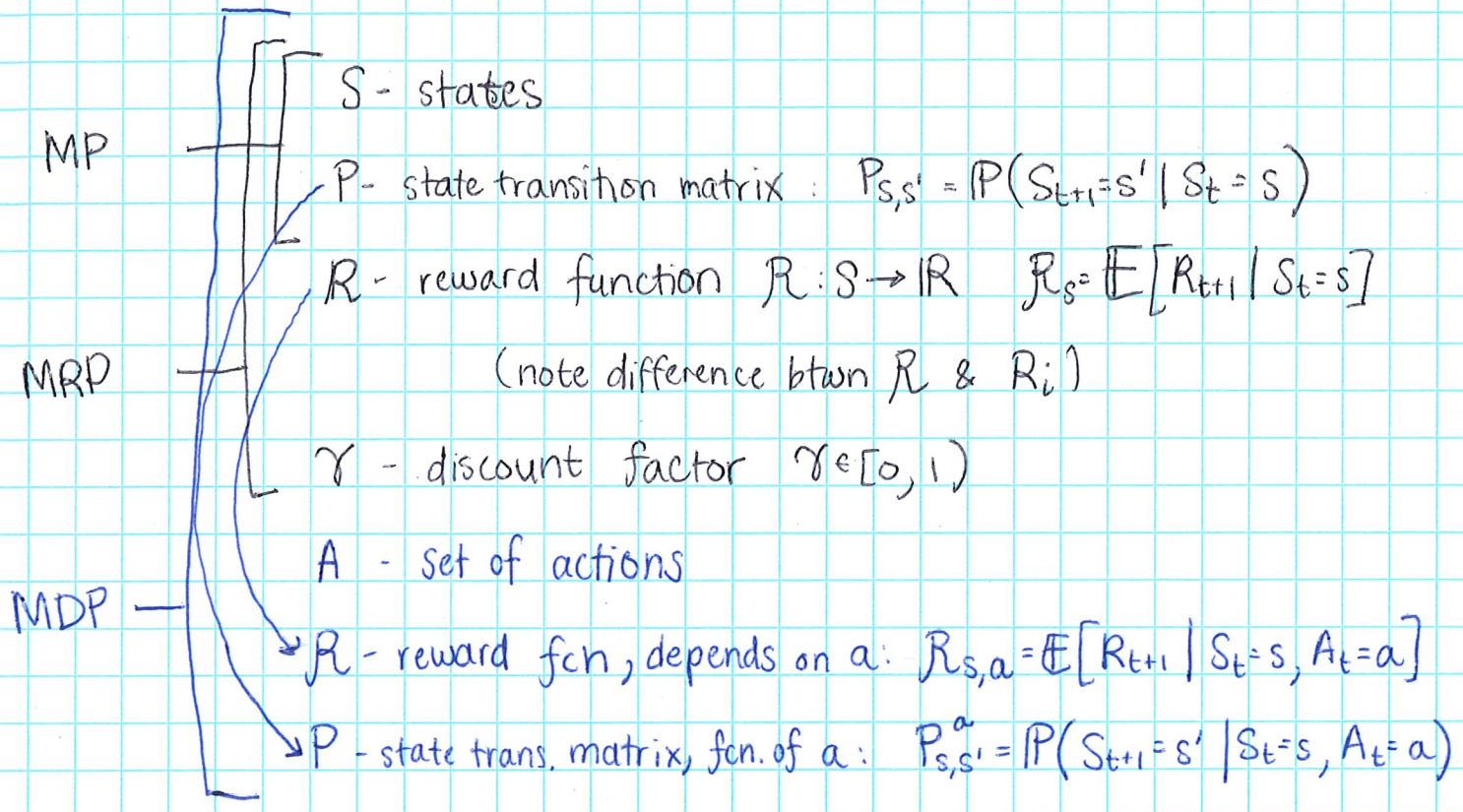


Markov Decision Processes



MP - Markov Process
 MRP - Markov Reward Process
 MDP - Markov Decision Process



MRP

State-value function: $v(s) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$

Expected discounted Rewards from s $\rightarrow \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s\right]$

$$\underline{v} = \begin{bmatrix} v(s_1) \\ \vdots \\ v(s_m) \end{bmatrix}$$

$$\underline{v} = \underline{R} + \gamma \cdot \underline{P} \cdot \underline{v}$$

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

algebra

MDP

Policy - $\pi(a|s) = P(A_t = a | S_t = s) \leftarrow \text{"Strategy"}$

For a fixed MDP & policy:

$\rightarrow \{S_i\}_{i=1}^{\infty}$ is a MP $\{(S_i, R_i)\}_{i=1}^{\infty}$ is a MRP

Ex: Given a policy π , ^{& MDP} how to compute transition matrix $P_{s,s'}$?

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

State-Value Function: $v_{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1}^{(\pi)} \mid S_t = s \right]$

$$= \mathbb{E}_{\pi} \left[R_{t+1}^{(\pi)} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s \right]$$

$$v_{\pi} = (\underline{I} - \gamma \cdot \underline{P}^{(\pi)})^{-1} \cdot \underline{R}^{(\pi)}$$

Action-Value Function:

$$q_{\pi}(s, a) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1}^{(\pi)} \mid S_t = s, A_t = a \right]$$

Q: How to optimize Rewards?

Bellman Equations

①

$$v^*(s) = \max_{\pi} v_{\pi}(s)$$

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

$$1) \quad v^*(s) = \max_a R_{s,a} + \gamma \cdot \sum_{s' \in S} P_{ss'}^a \cdot v^*(s')$$

$$q^*(s, a) = R_{s,a} + \gamma \sum_{s' \in S} P_{ss'}^a \max_{a'} q^*(s', a')$$

⇓

$$\pi^*(s) = \arg \max_a q^*(s, a)$$

Solve w/ iterative methods

- Value Iteration
- Q-learning