

history $\Rightarrow H_t = A_1, O_1, R_1, \dots, A_t, O_t, R_t$

State \rightarrow environment state S_t^e
 \rightarrow agent state S_t^a

* S_t is **Markov** $\Leftrightarrow P(S_{t+1} | S_t) = P(S_{t+1} | H_t)$

the state is a sufficient statistic of the future.

* environment types

① Fully Observable $O_t = S_t^a = S_t^e$
 (Markov Decision Process)

② Partially Observable (POMDP)

Method 1 $\rightarrow S_t^a = H_t$

naive approach

Method 2 $\rightarrow S_t^a = (P(S_t^e = s'), \dots, P(S_t^e = s''))$

bayesian probabilistic approach (beliefs)

Method 3 $\rightarrow S_t^a = \sigma(S_{t-1}^a W_s + O_t W_o)$

recurrent neural networks

an RL agent
can have

```
graph LR; A[an RL agent can have] --> B[Policy]; A --> C[Value function]; A --> D[Model];
```

policy

$$s \rightarrow \boxed{\pi} \rightarrow A$$

map from state to action

↳ deterministic $a = \pi(s)$

↳ probabilistic $\pi(a|s) = P[A=a | S=s]$

Value function

prediction of expected future reward

$$v_{\pi}(s) = E_{\pi} [R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots | S_t = s]$$

Model

predicts what environment will do next

↳ transition Model $P_{ss'}^a = P[S'=s' | S=s, A=a]$

↳ Rewards Model $R_s^a = E[R | S=s, A=a]$

Categorization of RL agents

① value-based (policy \times value \checkmark)

② policy-based (policy \checkmark value \times)

③ Actor-Critic (policy \checkmark value \checkmark)

① Model-free (policy/Value \checkmark model \times)

② Model-based (policy/Value \checkmark model \checkmark)

Sequential Decision Making

① reinforcement learning (through interaction)

↳ exploration vs exploitation

↳ prediction (given a policy) vs control (optimizing a policy)

② planning (no interaction needed)

Markov Decision Processes (MDP)

- ★ the environment is fully observable
- partially observable can be converted into MDP

Important Property of MDP: future is independent of the past, given the present

State transitions:

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

$$P = \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix} \quad \begin{array}{l} \text{State Transition Matrix} \\ \text{(each row sums to 1)} \end{array}$$

a Markov Process is a memoryless random Process
(a sequence of random states with the markov property)

$\langle S, P \rangle$

$\begin{array}{l} \downarrow \rightarrow \text{state transition probability matrix} \\ \rightarrow \text{(finite) state space} \end{array}$

Episode: a sample of the Markov Chain

Markov Reward Process

$\langle S, P, R, \gamma \rangle$

γ discount factor $\in [0, 1]$
reward function $R_s = E[R_{t+1} | S_t = s]$

Return G_t = total discounted reward at t

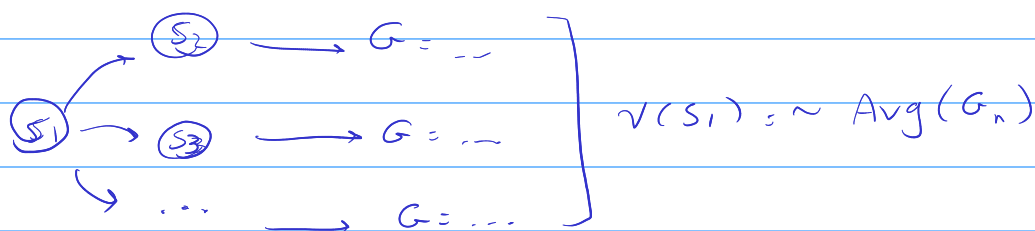
$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Myopic $0 \longleftarrow \gamma \longrightarrow 1$ far-sighted

value function : long-term value of state s

State-Value function : expected return starting from s

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



$$\begin{aligned} v(s) &= E[G_t | S_t = s] \\ &= E[R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= E[R_{t+1} + \gamma v(s_{t+1}) | S_t = s] \end{aligned}$$

Bellman Equation
for MDPs

$$v(s) = E[R_{t+1} + \gamma v(s_{t+1}) | S_t = s]$$

$$\rightarrow v(s) = R_s + \gamma \sum_{s' \in S} P_{ss'} v(s')$$

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

$$\begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix} = \begin{bmatrix} R_1 \\ \vdots \\ R_n \end{bmatrix} + \gamma \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & & \vdots \\ P_{m1} & \dots & P_{mn} \end{bmatrix} \begin{bmatrix} v(1) \\ \vdots \\ v(n) \end{bmatrix}$$

Matrix form of
Bellman Equation

now, Bellman equation is linear and solvable:

$$\underline{v = (I - \gamma P)^{-1} R} \quad O(n^3) \text{ for } n \text{ states}$$

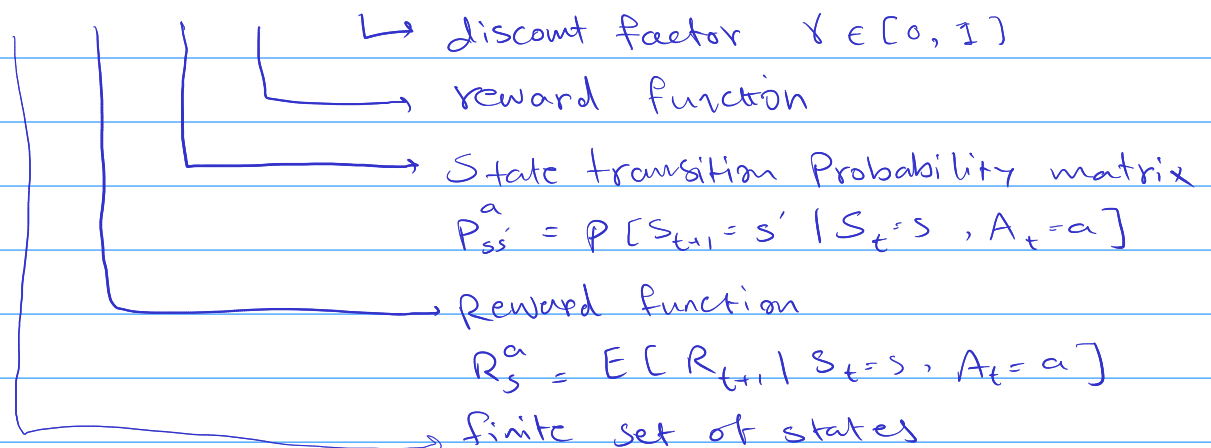
iterative methods

- ↳ Dynamic Programming
- ↳ Monte Carlo evaluation
- ↳ Temporal-Difference learning

Markov Decision Process : MDP with decisions

(still all states are markov)

$\langle S, A, P, R, \gamma \rangle$



Policy \rightarrow distributions over actions, given state

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

MDP \rightarrow state fully characterizes the path onwards

\hookrightarrow policies are stationary (time-independent)

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

$$S_1, R_2, S_2, \dots \Rightarrow \text{MRP}(\langle s, P^\pi, R, \gamma \rangle)$$

recovering MRP
from policy

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

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

state-value function $\rightarrow V_\pi(s) = E_\pi[G_t | S_t = s]$

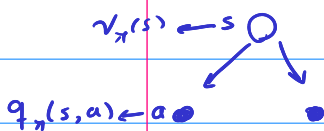
action-value function $\rightarrow Q_\pi(s, a) = E_\pi[G_t | S_t = s, A_t = a]$

if we're following a specific policy!

\Rightarrow Bellman Expectation Equation

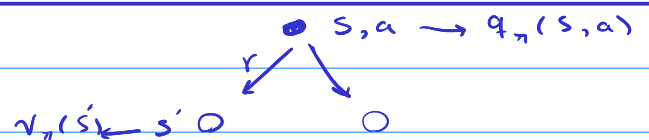
$$V_\pi(s) = E_\pi[R_{t+1} + \gamma V_\pi(S_{t+1}) | S_t = s]$$

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



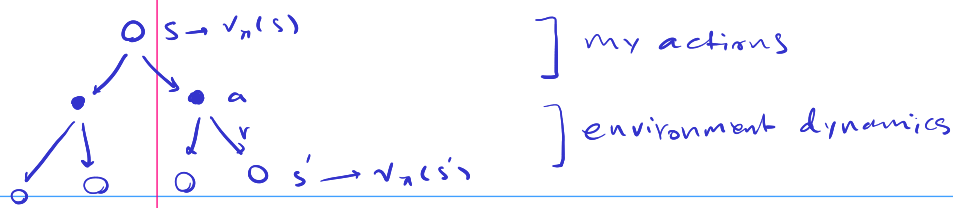
$$V_\pi(s) = \sum_{a \in A} \pi(a|s) Q_\pi(s, a)$$

calculate state value from action values by averaging over actions



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

average over dynamics of env



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

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

Optimal state-value func

Optimal action-value func

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

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

$$\pi \succ \pi' \text{ if } v_\pi(s) \geq v_{\pi'}(s) ; \forall s$$

* for any MDP there exists ^{deterministic} an optimal policy, π_* , $\pi: \forall \pi$

* there can exist multiple optimal policies

$$\pi_*(a|s) = \begin{cases} 1 & \text{if } a = \arg \max_{a \in A} q_*(s, a) \\ 0 & \text{otherwise} \end{cases}$$

$$v_*(s) = \max_a q_*(s, a)$$

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

$$v_*(s) = \max_a R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v_*(s')$$

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

Bellman Optimality Equations are non-linear

↳ value iteration

↳ policy iteration

↳ Q-learning

↳ sarsa