

"Solving" Finite MDPs

Solving: Finding policy that acts optimally in the MDP (with model free RL)

Finite horizon (episodic)

Infinite horizon (continuing)

Undiscounted*

$V(s)$ - Maximize expected future reward

for discrete states*

Discounted

Undiscounted

Discrete states

Continuous states

Tabular Methods*

Function Approximation

Value-based

Policy-based

Update V & Q using TD Learning
n-step return
TD(λ)

Update V & Q using Gradient Descent
n-step return
TD(λ)

Maximize expected future reward from start state*
REINFORCE
REINFORCE w/ baseline
Actor-Critic

Single step (Bandits)

Multi-step (Regular MDP)

Maximize expected total reward over time steps

Use action-value to select an action
 $q_*(a) = E[R_t | A_t = a]$
Choose action according to :
greedy, ϵ -greedy, UCB

Average Reward setting

Maximize average rate of reward / Average Reward

Stationary action distribution

Non-stationary action distribution

Sample average

$$Q_{n+1} = Q_n + \alpha(R_n - Q_n)$$
$$\alpha = 1/n$$

Exponential recency-weighted average

$$Q_{n+1} = (1 - \alpha)^n Q_1 + \sum_i \alpha (1 - \alpha)^{n-i} R_i$$
$$\alpha = \text{constant, in interval } (0, 1)$$