

Dynamic programming and Markov decision process

Consider a stationary MDP given by $(S, \mathcal{A}, P, R, \gamma)$

- * S is a finite set of states, $S_t \in S$ denotes the state at time t .
- * \mathcal{A} is a finite set of actions, $A_t \in \mathcal{A}$ is the action taken at time t .
- * P is a probability transition kernel: $P_{ss'}^a = P[S_{t+1} = s' | S_t = s, A_t = a]$
- * R is the reward function: $R_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$
- * γ is the discount factor involved in the return: $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$, $0 < \gamma < 1$

Optimal policies via Bellman optimality equations

value function $\rightarrow v^*(s) = \max_a (R_s^a + \gamma \sum_{s'} P_{ss'}^a v^*(s'))$

action value function $\rightarrow q^*(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a \max_{a'} q^*(s', a')$

Solving a MDP consists in finding v^* : this problem may be decomposed in a collection of subproblems, whose solutions are used many times to build the full solution (optimal substructure + overlapping problems).

\Rightarrow This justifies applying dynamic programming (DP).

Concretely: value iteration algorithm = iterative application of the Bellman optimal operator:

$$v_{k+1}(s) = (Tv_k)(s) = \max_a (R_s^a + \gamma \sum_{s'} P_{ss'}^a v_k(s'))$$

Naive implementation involve synchronous back ups, where all states are updated during each sweep. Asynchronous back ups can substantially improve computational cost. Example: in-place value iteration with prioritized ordering according to Bellman error: $|Tv(s) - s| = \epsilon(s)$.

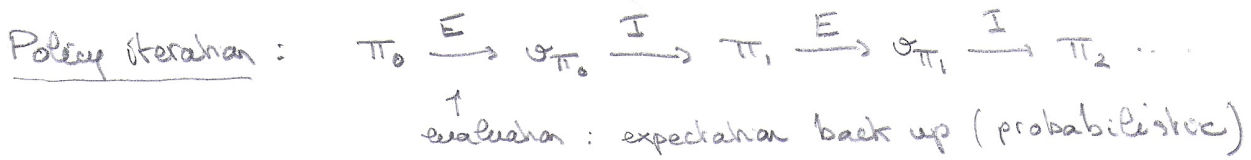
Key Limitation:

Dynamic programming methods assume knowledge of the environment dynamics. Monte Carlo methods offer to use sampling in order to improve policies without full knowledge of the environment. The central problem is then to balance exploitation and exploration.

Monte-Carlo method

improvement: greedy move (deterministic)

(2)



Sampling idea: The value function is the expected return conditioned on a starting state. Monte Carlo methods replace expectation by sampling means during the evaluation stage, assumed to last for a full episode (many samples)

Monte-Carlo and policy iteration coincides under assumption 1) of exploring starts (each action-state pair has non zero probability to be chosen at the start of an episode) and under assumption 2) of an infinite number of samples per episode.

Assumption 2) can be relaxed naturally as in value iteration.

Assumption 1) is relaxed via two distinct approaches: "on-policy" and off-policy approach.

On-policy approach: The same policy π is used for evaluation and control, but the policy is soft: $\min_a \pi(a|s) \geq \epsilon/|A|$

Example: ϵ -greedy policy π' select a random action uniformly with probability ϵ and the greedy action with probability $1-\epsilon$.

$$q_{\pi'}(s, \pi'(s)) = \sum_a \pi'(a|s) q_{\pi}(s, a) = \frac{\epsilon}{|A|} \sum_a q_{\pi}(s, a) + (1-\epsilon) \max_a q_{\pi}(s, a)$$

$$\geq \frac{\epsilon}{|A|} \sum_a q_{\pi}(s, a) + (1-\epsilon) \sum_a \frac{\pi(s, a) - \frac{\epsilon}{|A|}}{1-\epsilon} q_{\pi}(s, a) = v_{\pi}(s)$$

Off-policy approach: Distinct policies π and π' are used for evaluation and control. Typically π can be deterministic while π' ensures exploration.

number of visits to s in an episode \downarrow
 n_s
 $v_{\pi}(s) = \frac{1}{\sum_{i=1}^{n_s} p_i(s)}$
 probability of a sequence path after the i -th visit to s

$R_i(s) \leftarrow$ observed return from state s . only depends on policy taken

$$\frac{p_i(s)}{p'_i(s)} = \frac{\prod_{k=i}^{T-1} \pi(a_k|s_k) P_{s_k a_k s_{k+1}}^{ax}}{\prod_{k=i}^{T-1} \pi'(a_k|s_k) P_{s_k a_k s_{k+1}}^{ax}} = \frac{\prod_{k=i}^{T-1} \pi(a_k|s_k)}{\prod_{k=i}^{T-1} \pi'(a_k|s_k)}$$

Temporal difference algorithm

Monte Carlo methods differ from DP methods because:

- 1) they operate on sample experience: no need for environment model
- 2) they do not update their value estimates on the basis of other value estimates, i.e. they do not bootstrap (as in DP).

Temporal difference learning, perhaps the central original idea of reinforcement learning (v.s. control theory) combines Monte Carlo and DP approaches.

Central idea: bootstrapping should allow for online policy updates, e.g., avoid waiting for the end of an episode to perform evaluation. A simple, every visit Monte Carlo scheme for evaluation is

instantaneous reward used as target
= no bootstrapping →

$$v_{\pi}(s_t) \leftarrow v_{\pi}(s_t) + \alpha (R_{t+1}^{\pi} - v_{\pi}(s_t)), \quad \alpha: \text{fixed learning rate}$$

Alternative:

$$\begin{aligned} v_{\pi}(s) &= \mathbb{E}_{\pi} [G_t | S_t = s] \\ &= \mathbb{E}_{\pi} \left[\sum_{k \geq 0} \gamma^k R_{t+k+1} | S_t = s \right] \\ &= \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= \mathbb{E}_{\pi} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s] \end{aligned}$$

bootstrapping = using a target that depends on other estimates in the Monte-Carlo scheme

TD(0) rule:

$$v_{\pi}(s_t) \leftarrow v_{\pi}(s_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(s_t)]$$

Example of on-policy TD method: SARSA. For each episode

- 1) given current guess for $q(s, a)$, initialize s and choose a ϵ -greedily
- 2) observe r and s' , choose a' ϵ -greedily given s'
- 3) update $q(s, a) \leftarrow q(s, a) + \alpha [r + \gamma q(s', a') - q(s, a)]$.
- 4) $s \leftarrow s'$, $a \leftarrow a'$...

Q-learning = off-policy TD algorithm (Watkins, 1989) (4)

At the core of Q-learning is a simple value iteration update. That simplicity allows for simple convergence proof.

SARSA is on-policy because update can only be made on the basis of actions that are taken. Q-learning updates on the basis of the best available action possible.

For each episode:

1) initialize s

2) repeat for each step of the episode

a) choose a given s ϵ -greedily

b) take action a and observe s' and r

c) $q(s, a) \leftarrow q(s, a) + \alpha [r + \gamma \max_{a'} q(s', a') - q(s, a)]$

d) $s \leftarrow s'$ ↑ best possible action

3) terminate when halting condition is met.

Theorem: Given a finite MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$, the Q-learning algorithm given by the update rule:

$$q_{t+1}(s_t, a_t) = q_t(s_t, a_t) + \alpha_t(s_t, a_t) [r_{t+1} + \gamma \max_a q_t(s_{t+1}, a) - q_t(s_t, a_t)]$$

converges with probability 1 to the optimal action-value function as long as:

$$1) \sum_{t \geq 0} \alpha_t(s, a) = +\infty$$

$$2) \sum_{t \geq 0} \alpha_t^2(s, a) < +\infty$$

for all (s, a)

Sketch of the proof: Rewrite the update rule as

$$q_{t+1}(s_t, a_t) = (1 - \alpha_t) q_t(s_t, a_t) + \alpha_t [r_{t+1} + \gamma \max_a q_t(s_{t+1}, a)]$$

* Define $\Delta_t(s, a) = q_t(s, a) - q^*(s, a)$ where q^* is the optimal value function. Then the update rule reads

$$\Delta_{t+1}(s_t, a_t) = (1 - \alpha_t) \Delta_t(s_t, a_t) + \alpha_t [r_{t+1} + \gamma \max_a q(s_{t+1}, a) - q^*(s_t, a_t)]$$

* The goal is to show that Δ_t converges to zero with probability one.

The idea is the $F_t = r_{t+1} + \gamma \max_a q(s_{t+1}, a) - q^*(s_t, a_t)$ represents a random error that can be controlled. If the amplitude of the error are small enough, the relaxation term can confine Δ_t to zero. However α_t must not decrease to zero too fast for the relaxation to be strong enough.

* These ideas can be made precise in the context of stochastic approximation theory. Namely we have (Dvoretzky 1956)

The random process Δ_t taking value in \mathbb{R}^n and defined as

$$\Delta_{t+1}(s) = (1 - \alpha_t(s)) \Delta_t(s) + \alpha_t(s) F_t(s) \leftarrow \text{random noise}$$

converges to zero with probability one if:

$$1) \quad 0 \leq \alpha_t \leq 1 \quad \sum_{t \geq 0} \alpha_t(s) = +\infty \quad \text{and} \quad \sum_{t \geq 0} \alpha_t^2(s) < +\infty$$

$$2) \quad \mathbb{E}[F_t(s) | \mathcal{F}_t] \leq \gamma \|\Delta_t\|_\infty, \quad \gamma < 1$$

$$3) \quad \mathbb{V}[F_t(s) | \mathcal{F}_t] \leq C(1 + \|\Delta_t\|_\infty^2), \quad C \geq 0$$

* We only need to check 2) and 3). We will justify why 1) is a natural assumption heuristically.

First observe q^* is the fixed point of the contraction

$$\text{operator } (Tq)(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a \max_{a'} q(s', a')$$

Indeed, we have:

$$\begin{aligned} |Tq(s,a) - Tq'(s,a)| &= \gamma \left| \sum_{s'} P_{ss'}^a (\max_{a'} q(s',a') - \max_{a'} q'(s',a')) \right| \\ &\leq \gamma \sum_{s'} P_{ss'}^a \left| \max_{a'} q(s',a') - \max_{a'} q'(s',a') \right| \\ \text{triangular ineq.} \quad \rightarrow &\leq \gamma \sum_{s'} P_{ss'}^a \max_{a'} |q(s',a') - q'(s',a')| \\ &\leq \gamma \sum_{s'} P_{ss'}^a \|q - q'\|_{\infty} = \gamma \|q - q'\|_{\infty} \end{aligned}$$

Thus, we can evaluate:

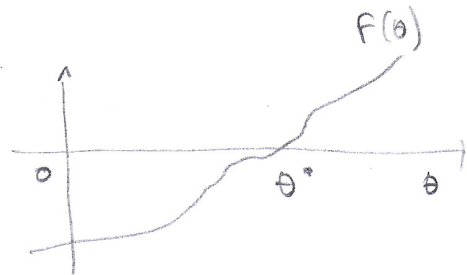
$$\begin{aligned} |E[F_b(s,a) | F_t]| &= |R_s^a + \gamma \sum_{s'} P_{ss'}^a \max_{a'} q(s',a') - q^+(s,a)| \\ &= |Tq(s,a) - Tq^+(s,a)| \leq \gamma \|q - q^+\|_{\infty} = \gamma \|\Delta_t\|_{\infty} \end{aligned}$$

This establishes 2). 3) immediately follows from boundedness of the reward, which we assume.

Heuristic justification of 1)

Consider a non-decreasing function:

Goal: finding a sequence of estimates
of the root θ^* when one has only
access to noisy values.



Without noise: Newton-Raphson $\theta_{n+1} = \theta_n - \frac{F(\theta_n)}{F'(\theta_n)}$
or $\theta_{n+1} = \theta_n - \alpha F(\theta_n)$ for α small enough

Idea: Can one come up with a sequence $\theta_n \rightarrow \theta^*$ by a
judicious choice of coefficient α_n such that

$$\theta_{n+1} = \theta_n - \alpha_n (F(\theta_n) - \gamma_n), \quad \gamma_n \text{ i.i.d. noise.}$$

Intuition: If f is linear, averaging the noise should be beneficial.

$$\text{This corresponds to taking } \alpha_n = \frac{1}{n+1} \Rightarrow \theta_n = \frac{\sum_{i=1}^n \gamma_i}{n}$$

If $\alpha_n = \frac{1}{n+1}$, one has $\sum_{n \geq 0} \alpha_n = +\infty$ and $\sum_{n \geq 0} \alpha_n^2 < +\infty$.

How general is this observation?

$$\begin{aligned}
 \hookrightarrow \text{consider } \Theta_{n+1} &= \Theta_n - \alpha_n (\Theta_n - \zeta_n) \\
 &= (1 - \alpha_n) \Theta_n + \alpha_n \zeta_n \\
 &= (1 - \alpha_n)(1 - \alpha_{n-1}) \Theta_{n-1} + \alpha_n \zeta_n + (1 - \alpha_n) \alpha_{n-1} \zeta_{n-1} \\
 &\vdots \\
 &= \underbrace{\prod_{i=1}^n (1 - \alpha_i)}_{\text{deterministic error}} \Theta_0 + \underbrace{\prod_{i=1}^n \prod_{k=i+1}^n (1 - \alpha_k)}_{\text{noise error}} \alpha_i \zeta_i
 \end{aligned}$$

for $\Theta_n \xrightarrow{L^2} \Theta^*$, we must have 1) $\lim_{n \rightarrow +\infty} \prod_{i=1}^n (1 - \alpha_i)^2 = 0$

$$2) \lim_{n \rightarrow +\infty} \sum_{i=1}^n \prod_{k=i+1}^n (1 - \alpha_k)^2 \alpha_i^2 = 0$$

1) is equivalent to $2 \sum_{i=1}^n \log(1 - \alpha_i) \sim -2 \sum_{i=1}^n \alpha_i \xrightarrow{n \rightarrow +\infty} -\infty$

2) follows from: $\exists N_\epsilon, n \geq N_\epsilon \Rightarrow \sum_{i=n}^{+\infty} \alpha_i^2 < \epsilon/2$

Thus taking α_i in $(0, 1)$, we have, for $n \geq N_\epsilon$

$$\begin{aligned}
 \left| \sum_{i=1}^n \prod_{k=i+1}^n (1 - \alpha_k)^2 \alpha_i^2 \right| &\leq \left| \sum_{i=1}^{N_\epsilon} \prod_{k=i+1}^n (1 - \alpha_k)^2 \alpha_i^2 \right| + \sum_{i=N_\epsilon}^{+\infty} \alpha_i^2 \\
 &\leq N_\epsilon \underbrace{\prod_{i=N_\epsilon}^n (1 - \alpha_i)^2}_{\xrightarrow{n \rightarrow +\infty} 0} + \epsilon/2 \leq \epsilon/2
 \end{aligned}$$