

facebook

Artificial Intelligence Research

How to solve an MDP incrementally: Reinforcement Learning

Alessandro Lazaric

Facebook AI Research

Outline

1 Policy Evaluation

2 Actor-Critic

3 Q-Learning

Policy Evaluation

Fixed policy π

For $i = 1, \dots, n$

1 Set $t = 0$

2 Set initial state s_0

3 While ($s_{t,i}$ not terminal) *[execute one trajectory]*

1 Take action $a_{t,i} = \pi(s_{t,i})$

2 Observe **next state** $s_{t+1,i}$ and **reward** $r_{t,i} = r(s_{t,i}, a_{t,i}) \sim \nu(s_{t,i}, a_{t,i})$

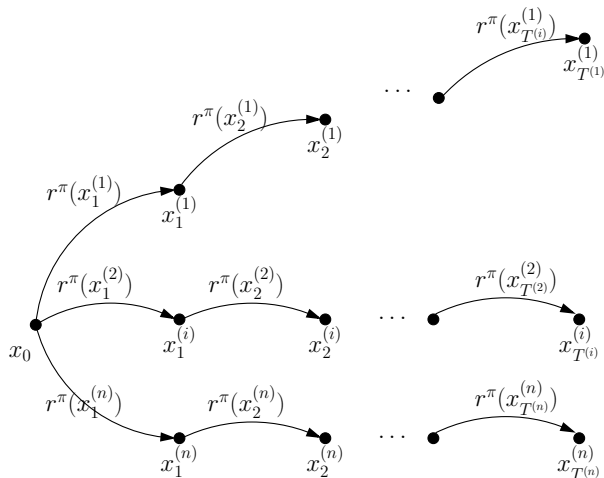
3 Set $t = t + 1$

EndWhile

EndFor

Return: Estimate of the value function $\hat{V}^\pi(\cdot)$

The RL Interaction Protocol



State Value Function

Fixed point of Bellman equation

$$V^\pi(s) = r^\pi(s) + \gamma \sum_{s'} p^\pi(s'|s) V^\pi(s')$$

Temporal Difference $TD(0)$

Input: π, T, s_0, V_0

$V = V_0, s = s_0$

for $t = 1, \dots, T$ **do**

 Execute action $a_t = \pi(s_t)$

 Observe r_t and next state s_{t+1}

 Update

$$\hat{V}(s_t) = \hat{V}(s_t) + \alpha(s_t)(r_t + \gamma\hat{V}(s_{t+1}) - \hat{V}(s_t))$$

end

return \hat{V}

Temporal Difference $TD(0)$: Intuition

The **Temporal difference error** of estimate \hat{V}^π w.r.t. transition (s_t, r_t, s_{t+1})

$$\delta_t = r_t + \gamma \hat{V}^\pi(s_{t+1}) - \hat{V}^\pi(s_t)$$

Temporal Difference $TD(0)$: Intuition

The **Temporal difference error** of estimate \hat{V}^π w.r.t. transition (s_t, r_t, s_{t+1})

$$\delta_t = r_t + \gamma \hat{V}^\pi(s_{t+1}) - \hat{V}^\pi(s_t)$$

- Bellman error for function \hat{V} at state s

$$\mathcal{B}^\pi(\hat{V}; s) = r^\pi(s) + \gamma \sum_{s'} p^\pi(s'|s) \hat{V}(s') - \hat{V}(s) \quad [\mathcal{B}^\pi(V^\pi; s) = 0]$$

- Conditioned on s_t , δ_t is an **unbiased** estimate of \mathcal{B}^π

$$\mathbb{E}_{r_t, s_{t+1}} [\delta_t | s_t] = r^\pi(s_t) + \gamma \mathbb{E}_{s_{t+1} | s_t} [\hat{V}^\pi(s_{t+1})] - \hat{V}^\pi(s_t) = \mathcal{B}^\pi(\hat{V}^\pi, s_t)$$

- Expected dynamics of TD(0)

$$\bar{V}^\pi(s_t) = \bar{V}^\pi(s_t) + \alpha(s_t) \mathcal{B}^\pi(\bar{V}^\pi; s_t)$$

Temporal-Difference TD(0): Alternative View

- Mix between old and new estimate of $V^\pi(s_t)$

old estimate $\hat{V}^\pi(s_t)$ new estimate $r_t + \gamma \hat{V}^\pi(s_{t+1})$

- Weighted average

$$\hat{V}^\pi(s_t) = (1 - \alpha) \hat{V}^\pi(s_t) + \alpha (r_t + \gamma \hat{V}^\pi(s_{t+1}))$$

Temporal-Difference TD(0): Properties

Theorem

Let TD(0) run with learning rate $\alpha(N_t(s_t))$ where $N_t(s_t)$ is the number of visits to the state s_t . If all states are visited *infinitely often* and the learning rate is set such that

$$\sum_{t=0}^{\infty} \alpha(t) = \infty \quad \sum_{t=0}^{\infty} \alpha(t)^2 < \infty \quad [\text{Robbins Monro's condition}]$$

then for any state $s \in \mathcal{S}$

$$\hat{V}^{\pi}(s) \xrightarrow{\text{a.s.}} V^{\pi}(s).$$

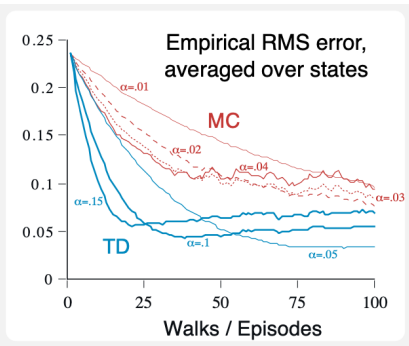
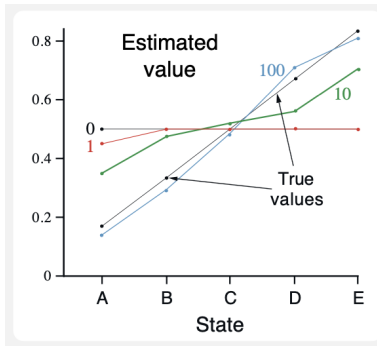
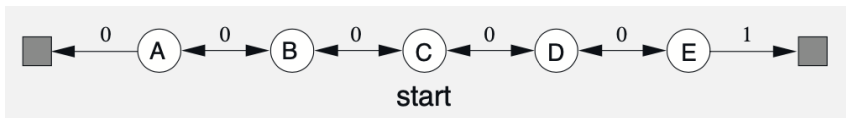
👉 standard choice is $\alpha_t(s) = \frac{1}{i^{\beta}}$ for $\beta \in (1/2, 1]$

$$\hat{V}_t(s) = \hat{V}_{t-1}(s) + \frac{1}{N_t(s)^{\beta}} (r + \gamma \hat{V}_{t-1}(s') - \hat{V}_{t-1}(s))$$

where $N_t(s)$ is the number of visits to state s before t

Example: Random Walk

(?)



Outline

1 Policy Evaluation

2 Actor-Critic

3 Q-Learning

Policy Learning

Learn optimal policy π^*

For $i = 1, \dots, n$

1 Set $t = 0$

2 Set initial state s_0

3 **While** ($s_{t,i}$ not terminal) *[execute one trajectory]*

1 **Take action** $a_{t,i}$

2 Observe next state $s_{t+1,i}$ and reward $r_{t,i} = r(s_{t,i}, a_{t,i})$

3 Set $t = t + 1$

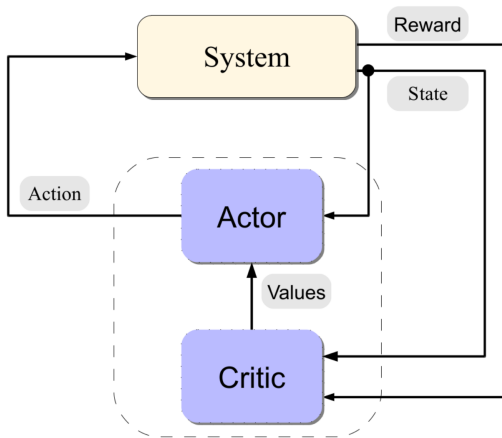
EndWhile

EndFor

Return: Estimate of the optimal policy $\hat{\pi}^*$

Actor-Critic Methods

- Critic: estimate the value function of the policy given by the actor
- Actor: change the policy to improve the value given by the critic



Policy Iteration

Let π_0 be an arbitrary stationary policy

while $k = 1, \dots, K$ **do**

Policy Evaluation: given π_k compute Q^{π_k}

Policy Improvement: find π_{k+1} that is better than π_k

 - e.g., compute the *greedy* policy

$$\pi_{k+1}(s) \in \arg \max_{a \in \mathcal{A}} Q^{\pi_k}(s, a)$$

return the last policy π_K

end

Policy Iteration is an instance of actor-critic methods:

- Actor: acts greedily (i.e., $\pi_{k+1} = \text{greedy}(Q^{\pi_k})$)
- Critic: compute Q given the policy π_k

Generalized Policy Iteration

More generally:

- Actor: *policy improvement* (even approximate)
- Critic: *policy evaluation*

(?) refers to actor/critic methods as **generalized policy iteration** (GPI)

👍 many possible schemes!

SARSA

Actor: the goal is to move the policy toward the greedy policy w.r.t. to the critic value

Possible strategies

- greedy (does not explore enough!)
- *ϵ -greedy* policy $\pi_Q(a|s)$

$$a = \mathcal{U}(\mathcal{A})$$

with probability ϵ

$$a = \arg \max_{\bar{a}} Q(s, \bar{a})$$

with probability $1 - \epsilon$

- *soft-max (random) exploratory* policy with temperature τ

$$\pi_Q(a|s) = \frac{\exp(Q(s, a)/\tau)}{\sum_{a'} \exp(Q(s, a')/\tau)}$$

The higher $Q(s, a)$, the more probability to take action a in state s

SARSA

Critic: uses TD to update the estimated Q-function

- Compute the temporal difference on the trajectory $\langle s_t, a_t, r_t, s_{t+1}, a_{t+1} \rangle$ (with actions chosen according to $\pi_Q(a|s)$)

$$\delta_t = r_t + \gamma \widehat{Q}(s_{t+1}, a_{t+1}) - \widehat{Q}(s_t, a_t)$$

SARSA

Critic: uses TD to update the estimated Q-function

- Compute the temporal difference on the trajectory $\langle s_t, a_t, r_t, s_{t+1}, a_{t+1} \rangle$ (with actions chosen according to $\pi_Q(a|s)$)

$$\delta_t = r_t + \gamma \hat{Q}(s_{t+1}, a_{t+1}) - \hat{Q}(s_t, a_t)$$

- Update the estimate of Q as

$$\hat{Q}(s_t, a_t) = \hat{Q}(s_t, a_t) + \alpha(s_t, a_t) \delta_t$$

State Action Reward State Action (SARSA) update

SARSA

Input: T, s_0, Q_0

$Q = Q_0, s = s_0, a = \mathcal{U}(\mathcal{A})$

for $t = 1, \dots, T$ **do**

 Execute action $a_t \sim \pi_{t,Q}(\cdot|s_t)$

 Observe r_t and next state s_{t+1}

$a_{t+1} \sim \pi_{t,Q}(\cdot|s_{t+1})$

 Update

$$\hat{Q}(s_t, a_t) = \hat{Q}(s_t, a_t) + \alpha(s_t)(r_t + \gamma \hat{Q}(s_{t+1}, a_{t+1}) - \hat{Q}(s_t, a_t))$$

end

return \hat{Q}, π_Q

Examples

$$\blacksquare \pi_{t,Q}(a|s) = \frac{\exp(Q(s, a)/\tau_t)}{\sum_b \exp(Q(s, b)/\tau_t)}$$

$$\blacksquare \pi_{t,Q}(a|s) = (1 - \epsilon_t) \mathbb{1} \left(a_t = \arg \max_b Q(s, b) \right) + \epsilon_t / A$$

SARSA: Properties (informal)

- The TD updates make \hat{Q} converge to Q^π
- The update of π_Q allows to improve the policy
- A decreasing temperature allows to become more and more greedy
 \Rightarrow If $\tau \rightarrow 0$ with a proper rate, then $\hat{Q} \rightarrow Q^*$ and $\pi_Q \rightarrow \pi^*$
- Similarly for ϵ -greedy

SARSA: Limitations

The actions a_t need to be selected according to the current Q

⇒ **On-policy learning**

Outline

1 Policy Evaluation

2 Actor-Critic

3 Q-Learning

The Optimal Bellman Equation

Proposition

The optimal value function Q^* (i.e., $Q^* = \max_{\pi} Q^{\pi}$) is the solution to the *optimal Bellman equation*:

$$Q^*(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|s, a) \max_{a' \in A} Q^*(s', a').$$

Q-Learning

Input: T, s_0, Q_0

$Q = Q_0, s = s_0$

for $t = 1, \dots, T$ **do**

 Execute action $a_t \sim \pi_t(s_t)$

 Observe r_t and next state s_{t+1}

 Update

$$\hat{Q}(s_t, a_t) = \hat{Q}(s_t, a_t) + \alpha(s_t, a_t)(r_t + \gamma \max_{a'} \hat{Q}(s_{t+1}, a') - \hat{Q}(s_t, a_t))$$

end

return $\hat{Q}, \text{greedy}(\hat{Q})$

The policy π_t is only required to cover all actions with non-zero probability

Q-Learning: Properties

Proposition

If the learning rate satisfies the Robbins-Monro conditions in all states $s, a \in S \times A$

$$\sum_{i=0}^{\infty} \alpha_i(s, a) = \infty, \quad \sum_{i=0}^{\infty} \alpha_i^2(s, a) < \infty,$$

and all state-action pairs are tried *infinitely often*, then for all $s, a \in S \times A$

$$\hat{Q}(s, a) \xrightarrow{a.s.} Q^*(s, a)$$

Remark: this is another example of R-M algorithm

Remark: “infinitely often” requires a steady exploration policy

Q-Learning vs. SARSA

Update rule

- *Q-Learning*

$$\widehat{Q}(s_t, a_t) = \widehat{Q}(s_t, a_t) + \alpha(r_t + \gamma \max_{a'} \widehat{Q}(s_{t+1}, a') - \widehat{Q}(s_t, a_t))$$

- *SARSA*

$$\widehat{Q}(s_t, a_t) = \widehat{Q}(s_t, a_t) + \alpha(r_t + \gamma \widehat{Q}(s_{t+1}, a_{t+1}) - \widehat{Q}(s_t, a_t))$$

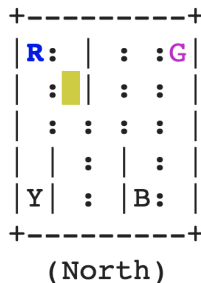
Both aim at learning the target policy $\pi^*(s) = \arg \max_a Q^*(s, a)$

- Q-learning can use any behavioral policy (*off-policy learning*)
- for SARSA the behavioral policy should be close to the estimated target policy (*on-policy learning*)

The Taxi Problem

Hands-on session [link] – [20min]

- Plot learning progress during training and evaluate performance at testing over multiple episodes
- How do learning rate and exploration parameters influence the performance?





Thank you!

facebook

Artificial Intelligence Research



. \ |