

Policy Iteration on Linear Quadratic Problem



Farnaz Adib Yaghmaie

Linköping University, Sweden
farnaz.adib.yaghmaie@liu.se

April 6, 2021

■ Dynamics:

$$s_{t+1} = As_t + Bu_t + w_t$$

■ State and action:

$$s_t \in \mathbb{R}^n,$$

$$u_t \in \mathbb{R}^m$$

■ Cost function (\equiv negative of reward):

$$c_t = s_t^\top Q s_t + u_t^\top R u_t, \quad Q \geq 0, R > 0$$

Solvability Criterion: Minimize $V(s)$ with respect to the policy π

$$V(s_t) = \mathbf{E}\left[\sum_{k=t}^{+\infty} (c_k - \lambda) | s_t\right]$$

where λ is the average cost

$$\lambda = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T c_t$$

Note that minimizing $V(s_t)$ and λ are equivalent.

The agents learn a quadratic Q function

$$Q(s, a) = \begin{bmatrix} s^\dagger & a^\dagger \end{bmatrix} \begin{bmatrix} g_{ss} & g_{sa} \\ g_{sa}^\dagger & g_{aa} \end{bmatrix} \begin{bmatrix} s \\ a \end{bmatrix} = z^\dagger G z \quad (1)$$

The policy is given by optimizing the Q function

$$\pi = -g_{aa}^{-1} g_{sa}^\dagger s = K s \quad (2)$$

- 1 Compute the empirical average cost $\lambda = \frac{1}{T} \sum_{t=1}^T c_t$
- 2 Collect data
 - Observe s and select a

$$a = Ks + r, \quad r \sim \mathcal{N}(0, \sigma^2).$$

- Apply a and observe r, s' .
 - Add s, a, r, s' to the history.
- 3 Estimated the kernel of Q by Least Squares Temporal Difference (LSTD)

$$\text{vecs}(G) = \left(\frac{1}{T} \sum_{t=1}^T \psi_t (\psi_t - \psi_{t+1})^\dagger \right)^{-1} \left(\frac{1}{T} \sum_{t=1}^T \psi_t (c_t - \lambda) \right) \quad (3)$$

where

$$z = \begin{bmatrix} s \\ a \end{bmatrix}, \quad \Psi = [z_1^2, 2z_1z_2, \dots, 2z_1z_n, z_2^2, \dots, 2z_2z_n, \dots, z_n^2]^\dagger.$$

■ Dynamics:

$$s_{t+1} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} s_t + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u_t + w_t$$

■ Cost function (\equiv negative of reward):

$$c_t = s_t^\top \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} s_t + u_t^\top 1 u_t.$$

- **Exact analytical solution assuming full information about dynamics**

$$u_t^* = \begin{bmatrix} -0.422 & -1.244 \end{bmatrix} s_t$$

- **Initialization of the algorithm**

$$u_t = \begin{bmatrix} -0.616 & -1.614 \end{bmatrix} s_t$$

Try the following:

- Run

`Crash_course_on_RL/q_on_lq_notebook.ipynb`

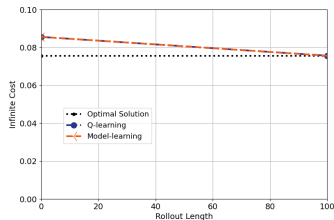
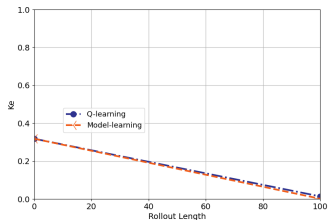
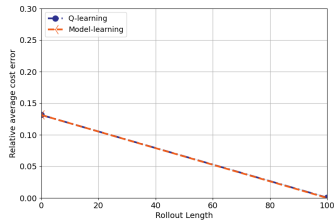
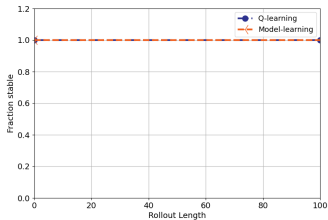
and verify the median of the error in K is $\sim 0.01\%$.

- Set

`'explore_mag=0.0'` in `'My_q_learning.ql'`

and verify that the agent cannot solve the problem! you don't get any stable controller.

- Make sure you understand the code!



Important observations

- Q-learning performs superb on the LQ problem
- No hyper-parameters to tune

Email your questions to

farnaz.adib.yaghmaie@liu.se