

Robot Learning

Winter Semester 2020/2021, Homework 2

Prof. Dr. J. Peters, J. Watson, J. Carvalho, J. Urain and T. Dam



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Total points: 20

Due date: Midnight, Monday, 7 December, 2020

Name, Surname, ID Number

Problem 2.1 Optimal Control [20 Points]

In this exercise, we consider a finite-horizon discrete time-varying Stochastic Linear Quadratic Regulator with Gaussian noise and time-varying quadratic reward function. Such system is defined as

$$\mathbf{s}_{t+1} = \mathbf{A}_t \mathbf{s}_t + \mathbf{B}_t \mathbf{a}_t + \mathbf{w}_t, \quad (1)$$

where \mathbf{s}_t is the state, \mathbf{a}_t is the control signal, $\mathbf{w}_t \sim \mathcal{N}(\mathbf{b}_t, \Sigma_t)$ is Gaussian additive noise with mean \mathbf{b}_t and covariance Σ_t and $t = 0, 1, \dots, T$ is the time horizon. The control signal \mathbf{a}_t is computed as

$$\mathbf{a}_t = -\mathbf{K}_t \mathbf{s}_t + \mathbf{k}_t \quad (2)$$

and the reward function is

$$\text{reward}_t = \begin{cases} -(\mathbf{s}_t - \mathbf{r}_t)^\top \mathbf{R}_t (\mathbf{s}_t - \mathbf{r}_t) - \mathbf{a}_t^\top \mathbf{H}_t \mathbf{a}_t & \text{when } t = 0, 1, \dots, T-1 \\ -(\mathbf{s}_t - \mathbf{r}_t)^\top \mathbf{R}_t (\mathbf{s}_t - \mathbf{r}_t) & \text{when } t = T \end{cases} \quad (3)$$

a) Implementation [8 Points]

Implement the LQR with the following properties

$$\begin{aligned} \mathbf{s}_0 &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}) & T &= 50 \\ \mathbf{A}_t &= \begin{bmatrix} 1 & 0.1 \\ 0 & 1 \end{bmatrix} & \mathbf{B}_t &= \begin{bmatrix} 0 \\ 0.1 \end{bmatrix} \\ \mathbf{b}_t &= \begin{bmatrix} 5 \\ 0 \end{bmatrix} & \Sigma_t &= \begin{bmatrix} 0.01 & 0 \\ 0 & 0.01 \end{bmatrix} \\ \mathbf{K}_t &= \begin{bmatrix} 5 & 0.3 \end{bmatrix} & \mathbf{k}_t &= 0.3 \\ \mathbf{H}_t &= 1 & \mathbf{R}_t &= \begin{cases} \begin{bmatrix} 100000 & 0 \\ 0 & 0.1 \end{bmatrix} & \text{if } t = 14 \text{ or } 40 \\ \begin{bmatrix} 0.01 & 0 \\ 0 & 0.1 \end{bmatrix} & \text{otherwise} \end{cases} & \mathbf{r}_t &= \begin{cases} \begin{bmatrix} 10 \\ 0 \end{bmatrix} & \text{if } t = 0, 1, \dots, 14 \\ \begin{bmatrix} 20 \\ 0 \end{bmatrix} & \text{if } t = 15, 16, \dots, T \end{cases} \end{aligned}$$

Execute the system 20 times. Plot the mean and 95% confidence (see “68–95–99.7 rule” and matplotlib.pyplot.fill_between function) over the different experiments of the state \mathbf{s}_t and of the control signal \mathbf{a}_t over time. How does the system behave? Compute and write down the mean and the standard deviation of the cumulative reward over the experiments. Attach a snippet of your code.

b) LQR as a P controller [4 Points]

The LQR can also be seen as a simple P controller of the form

$$a_t = \mathbf{K}_t (\mathbf{s}_t^{\text{des}} - \mathbf{s}_t) + k_t, \quad (4)$$

which corresponds to the controller used in the canonical LQR system with the introduction of the target $\mathbf{s}_t^{\text{des}}$. Assume as target

$$\mathbf{s}_t^{\text{des}} = \mathbf{r}_t = \begin{cases} \begin{bmatrix} 10 \\ 0 \end{bmatrix} & \text{if } t = 0, 1, \dots, 14 \\ \begin{bmatrix} 20 \\ 0 \end{bmatrix} & \text{if } t = 15, 16, \dots, T \end{cases} \quad (5)$$

Use the same LQR system as in the previous exercise and run 20 experiments. Plot in one figure the mean and 95% confidence (see “68–95–99.7 rule” and `matplotlib.pyplot.fill_between` function) of the first dimension of the state, for both $\mathbf{s}_t^{\text{des}} = \mathbf{r}_t$ and $\mathbf{s}_t^{\text{des}} = \mathbf{0}$.

c) Optimal LQR [8 Points]

To compute the optimal gains \mathbf{K}_t and \mathbf{k}_t , which maximize the cumulative reward, we can use an analytic optimal solution. This controller recursively computes the optimal action by

$$\mathbf{a}_t^* = -(\mathbf{H}_t + \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T (\mathbf{V}_{t+1} (\mathbf{A}_t \mathbf{s}_t + \mathbf{b}_t) - v_{t+1}), \quad (6)$$

which can be decomposed into

$$\mathbf{K}_t = -(\mathbf{H}_t + \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{A}_t, \quad (7)$$

$$\mathbf{k}_t = -(\mathbf{H}_t + \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T (\mathbf{V}_{t+1} \mathbf{b}_t - v_{t+1}). \quad (8)$$

where

$$\mathbf{M}_t = \mathbf{B}_t (\mathbf{H}_t + \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{B}_t)^{-1} \mathbf{B}_t^T \mathbf{V}_{t+1} \mathbf{A}_t \quad (9)$$

$$\mathbf{V}_t = \begin{cases} \mathbf{R}_t + (\mathbf{A}_t - \mathbf{M}_t)^T \mathbf{V}_{t+1} \mathbf{A}_t & \text{when } t = 1 \dots T-1 \\ \mathbf{R}_t & \text{when } t = T \end{cases} \quad (10)$$

$$v_t = \begin{cases} \mathbf{R}_t \mathbf{r}_t + (\mathbf{A}_t - \mathbf{M}_t)^T (\mathbf{V}_{t+1} \mathbf{b}_t - v_{t+1}) & \text{when } t = 1 \dots T-1 \\ \mathbf{R}_t \mathbf{r}_t & \text{when } t = T \end{cases} \quad (11)$$

Run 20 experiments with: the controller from a); the P controller from b) with $\mathbf{s}_t^{\text{des}} = \mathbf{r}_t$ (with \mathbf{r}_t as defined in a)); and with the controller c) resulting from computing the optimal gains \mathbf{K}_t and \mathbf{k}_t . Plot the mean and 95% confidence (see “68–95–99.7 rule” and `matplotlib.pyplot.fill_between` function) of both states for all three different controllers used so far. Use one figure per state. Report the mean and std of the cumulative reward for each controller and comment the results. Attach a snippet of your code.