

Robot Learning

Group 80: Lukas Schneider(2565695), Si Jun Kwon(2453134), Michael Erni(2547105)



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Sommersemester 2020
Sheet 2

Task 2.1

2.1a)

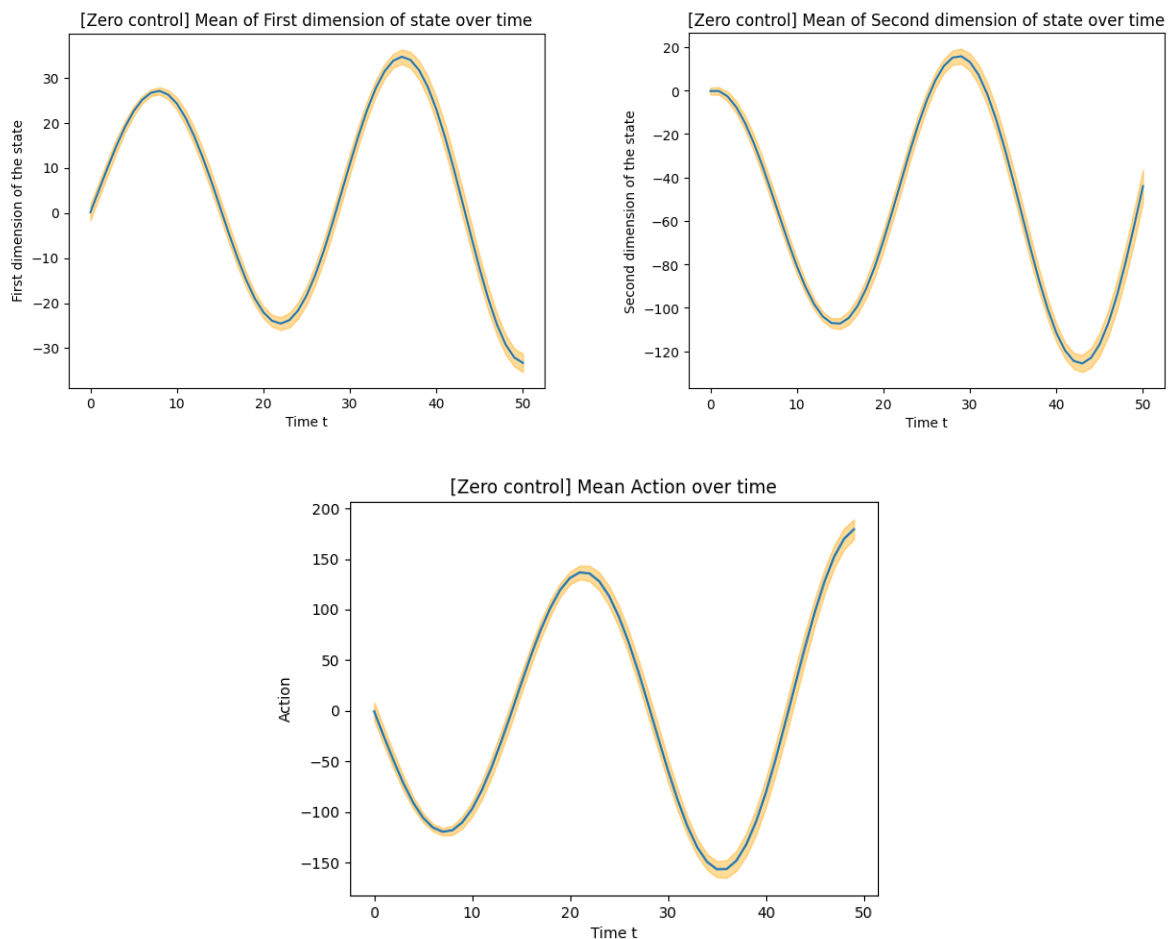


Figure 1: Mean position and 95% confidence interval of the joints with controller a), mean action of the controller a)

The system described in this subtask can be seen as a P controller with $s_t^{\text{des}} = 0$ for all $t = 0, 1, \dots, T$. In fact if we take the formula from b and set $s_t^{\text{des}} = 0$ for all $t = 0, 1, \dots, T$ we get the same formula to calculate the action in a). Thus the controller in this subtask does not consider the desired state r_t but instead desires for the point $(0 \ 0)^T$. It is also

apparent from the plots that the gains are too great and the system is instable hinting at suboptimal gains K_T and offset k_t . First the system drifts off in the positive direction of the first dimension because of the drift component b_t of the modeled SLQR. To compensate this the system takes action into the negative direction (refer to Figure ??). This propells the second dimension of state into the negative direction (at first the first dimension stays unaffected from the actions as the actions don't affect the first dimension directly). Due to the influence of the second joint on the first joint (refer to $A_{t,(1,2)} = 0.1$) the first joint slowly accelerates into negative direction. It then overshoots $s_t^{\text{des}} = 0$ because the second joint does not consider the velocity approaching s_t^{des} . Analogously the same happens for the compensation in the positive direction. This sequence of events repeats itself, the system becoming more instable with each repetition (the maximal distance to $(0\ 0)^T$ becoming bigger).

The average cumulative reward for this system was -2976946.972171863 with a std of 827060.1944006783.

Code

```

6 class LQRSim(abc.ABC):
7     """
8     Class that simulate the LQR system
9     """
10    def __init__(self) -> None:
11        self.s_hist = np.array(
12            [np.random.multivariate_normal(np.zeros((2,)), np.eye(2))]
13        ).reshape(2, 1)
14        self.a_hist = np.zeros((0, ))
15        self.T = 50
16        self.A_t = np.array([
17            [1, .1],
18            [0, 1]
19        ])
20        self.B_t = np.array([0, .1]).reshape(-1, 1)
21        self.b_t = np.array([5., 0.])
22        self.Sigma_T = np.array([
23            [.01, 0],
24            [0, .01]
25        ])
26        self.K_t = np.array([5, .3]).reshape(1, -1)
27        self.k_t = .3
28        self.H_t = 1
29        self.R_t = np.array([
30            [.01, 0],
31            [0, .1]
32        ]) * self.T
33        self.R_t[14] = np.array([
34            [100000, 0],
35            [0, .1]
36        ])
37        self.R_t[40] = np.array([
38            [100000, 0],
39            [0, .1]
40        ])
41        self.r_t = np.append(
42            np.array([[10.], [0.]] * 15),
43            np.array([[20.], [0.]] * 36),
44            axis=0
45        ).T[0]
46        self.ran = False
47
48    @abc.abstractclassmethod
49    def step_func(self, t):
50        pass
51
52    def run(self):
53        """
54        runs the simulation
55        """
56        if self.ran:
57            return

```

```

58     for t in range(self.T):
59         current_state = self.s_hist[:, t].reshape(-1, 1)
60         action = self.step_func(t)
61         noise = np.random.multivariate_normal(self.b_t, self.Sigma_T)\
62             .reshape(-1, 1)
63         next_state = self.A_t@current_state + self.B_t*action + noise
64         self.s_hist = np.append(self.s_hist, next_state, axis=1)
65         self.a_hist = np.append(self.a_hist, action)
66     self.ran = True
67
68     def calc_reward(self):
69         if not self.ran:
70             raise Exception(
71                 "The simulation must be run before rewards can calculated"
72             )
73         self.reward = np.zeros((self.T, ))
74         err = (self.s_hist[:, -1] - self.r_t[:, -1]).reshape(-1, 1)
75         self.reward[-1] = -err.T@self.R_t[-1]@err
76         for t in range(self.T - 2, -1, -1):
77             err = (self.s_hist[:, t] - self.r_t[:, t])
78             self.reward[t] = -err.T@self.R_t[t]@err - self.a_hist[t].T * \
79                 self.H_t*self.a_hist[t]
80         return self.reward

```

Listing 1: Base SLQR class that simulates the system

This abstract class summarizes all the common variables and functions used by all tasks commonly. These consists of the static system matrices like A_t , the calculation of the next state given a action taken and calculation of the reward. The subclasses which derive from this class would implement the stepping function that return the action. The stepping function for the controller a is shown in the code listing below.

```

83 class LQR1(LQRSim):
84     def step_func(self, t):
85         current_state = self.s_hist[:, t].reshape(-1, 1)
86         action = -self.K_t@current_state + self.k_t
87         return action

```

Listing 2: SLQR with controller a

Finally a plotter is used plotting the mean position of the joints, mean action and the mean cumulative rewards. Refer to the listing below.

```

131 class Plotter:
132     i = 0
133
134     def plot(self, sims, s1=True, s2=True, a=True, r=True, prefix=''):
135         to_plot = {}
136         if s1:
137             to_plot['s1'] = self.i
138             self.i += 1
139
140         if s2:
141             to_plot['s2'] = self.i
142             self.i += 1
143
144         if a:
145             to_plot['a'] = self.i
146             self.i += 1
147
148         if r:
149             to_plot['r'] = self.i
150             self.i += 1
151
152         time_states = np.linspace(0, 50, 51)
153         time_r_a = np.linspace(0, 49, 50)
154         [plt.figure(ii) for ii in range(self.i)]
155
156         states = np.array([sim.s_hist for sim in sims])

```

```

157 states_mean = np.mean(states , axis=0)
158 states_var = np.var(states , axis=0)
159
160 actions = np.array([sim.a_hist for sim in sims])
161 action_mean = np.mean(actions , axis=0)
162 action_var = np.var(actions , axis=0)
163
164 cumulative_rewards = np.array(
165     [np.cumsum(sim.calc_reward()) for sim in sims]
166 )
167 reward_mean = np.mean(cumulative_rewards , axis=0)
168 reward_std = np.sqrt(np.var(cumulative_rewards , axis=0))
169 print(f'{{prefix}} Mean cumulative reward {reward_mean[-1]}')
170 print(f'{{prefix}} Variance of cumulative reward {reward_std[-1]}')
171
172 if s1:
173     plt.figure(to_plot['s1'])
174     plt.title(f'{{prefix}} Mean of First dimension of state over time')
175     plt.plot(time_states , states_mean[0])
176     plt.fill_between(
177         time_states ,
178         states_mean[0] - 2 * np.sqrt(states_var[0]),
179         states_mean[0] + 2 * np.sqrt(states_var[0]),
180         color='orange',
181         alpha=.4
182     )
183     plt.xlabel('Time t')
184     plt.ylabel('First dimension of the state')
185
186 if s2:
187     plt.figure(to_plot['s2'])
188     plt.title(
189         f'{{prefix}} Mean of Second dimension of state over time'
190     )
191     plt.plot(time_states , states_mean[1])
192     plt.fill_between(
193         time_states ,
194         states_mean[1] - 2 * np.sqrt(states_var[1]),
195         states_mean[1] + 2 * np.sqrt(states_var[1]),
196         color='orange',
197         alpha=.4
198     )
199     plt.xlabel('Time t')
200     plt.ylabel('Second dimension of the state')
201
202 if a:
203     plt.figure(to_plot['a'])
204     plt.title(f'{{prefix}} Mean Action over time')
205     plt.plot(time_r_a , action_mean)
206     plt.fill_between(
207         time_r_a ,
208         action_mean - 2 * np.sqrt(action_var),
209         action_mean + 2 * np.sqrt(action_var),
210         color='orange',
211         alpha=.4
212     )
213     plt.xlabel('Time t')
214     plt.ylabel('Action')
215
216 if r:
217     plt.figure(to_plot['r'])
218     plt.title(f'{{prefix}} Mean Cumulative Reward over time')
219     plt.plot(time_r_a , reward_mean)
220     plt.xlabel('Time t')
221     plt.ylabel('Reward')

```

Listing 3: Plotter for the simulations

2.1b)

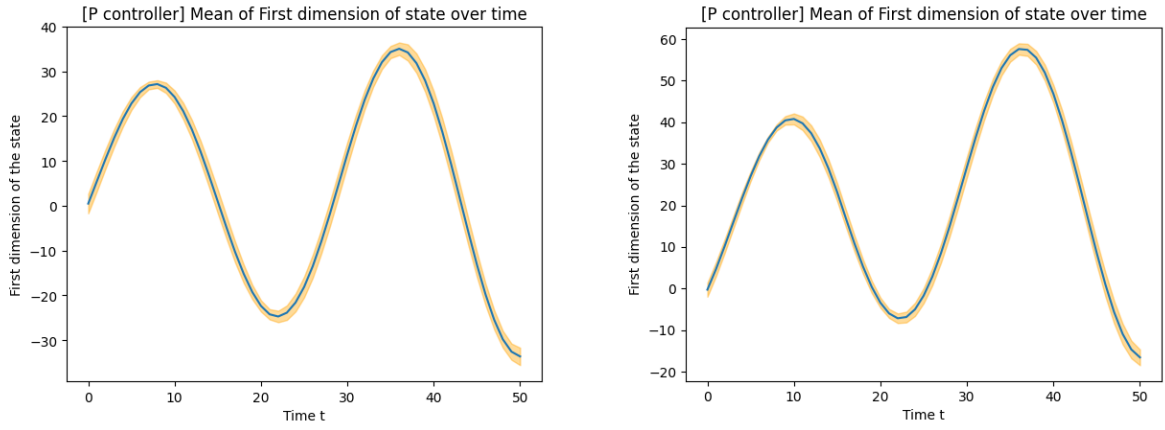


Figure 2: P controller with $s_t^{\text{des}} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ (left) and with $s_t^{\text{des}} = r_t$ (right)

2.1c)

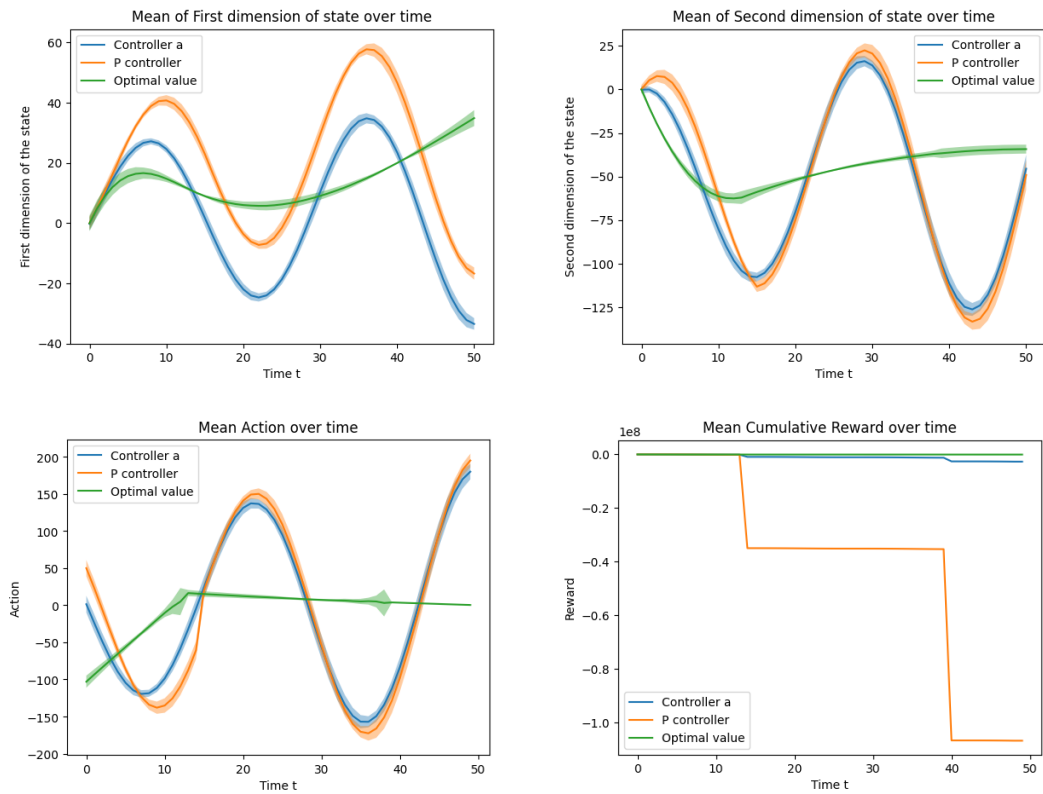


Figure 3: All three controllers

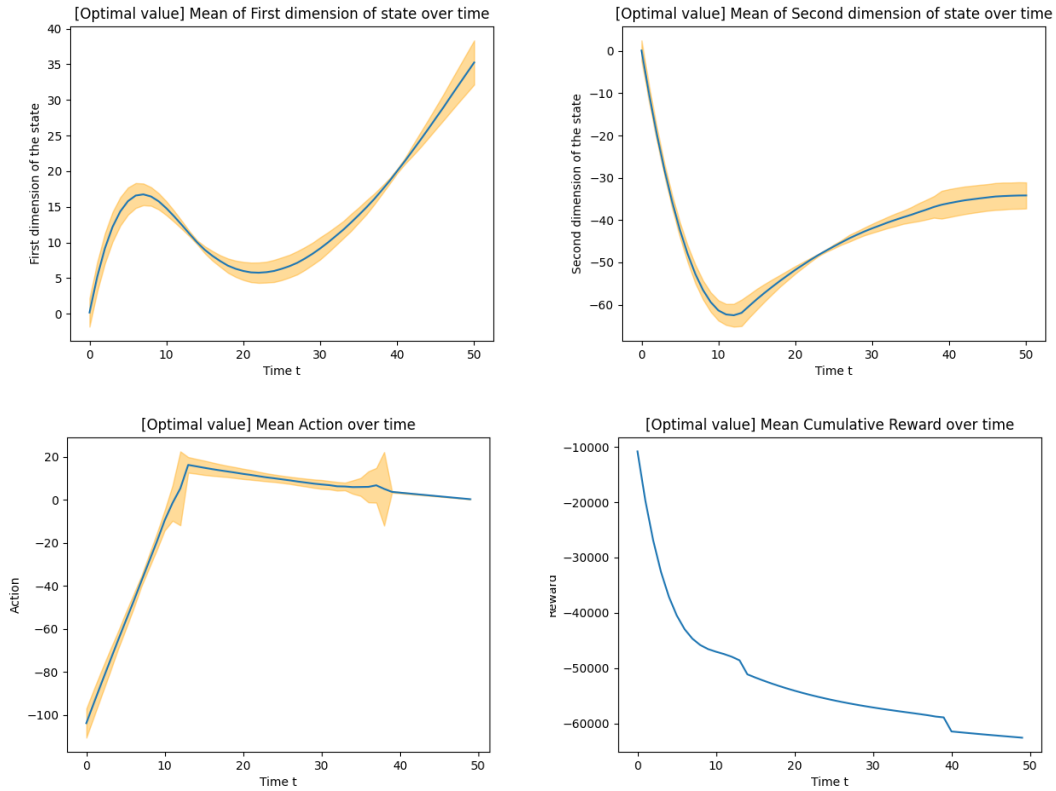


Figure 4: Optimal value controller

Controller type	Mean reward(the bigger the better)	Standard deviation of the rewards
Optimal Controller	-62584.876481951775	5182.129498148187
Controller a/P Controller($s_t^{\text{des}} = 0$)	-2976946.972171863	827060.1944006783
P Controller($s_t^{\text{des}} = r_t$)	-105344044.62913443	8505038.421848932

Figure 5: Table of mean rewards and standard deviation of the rewards in descending order of mean rewards

Looking at each controller one can make following observations:

- The Optimal controller outperforms the other variants by an outstanding amount. This is due to the fact that the optimal controller does not just take the current error into account but also the long-term reward and maximizes the long-term reward. This is shown by the trajectory of the system being controlled by the optimal controller. The controller overshoots by a little margin first so it can precisely hit the desired position (for the first joint as the first joint gives the biggest reward) at the time point where the reward is maximal($t = 14, 40$). Furthermore the optimal controller delivers is more consistent with its rewards as the standard deviation of rewards is fairly small. Refer to Figure 3 and Figure 4.
- Controller a) does not consider the desired trajectory but still manages to hit close to the desired point at the maximal reward time points($t = 14, 40$) by sheer coincidence(it desires for the point $\begin{pmatrix} 0 \\ 0 \end{pmatrix}$). Refer to Figure 3.
- Ironically the P controller with $s_t^{\text{des}} = r_t$ performs worse than the Controller a). This is due to the overshooting problem of the P controller. This imprecision is amplified with the given desired trajectory. Refer to Figure 3.

Refer to the next page for the codes used for this subtask.

Code

```
90 class LQR_PCONTROLLER(LQRSim):
91     def __init__(self, s_des=None) -> None:
92         super().__init__()
93         self.s_des = s_des if s_des is not None else self.r_t
94
95     def step_func(self, t):
96         current_state = self.s_hist[:, t].reshape(-1, 1)
97         action = self.K_t@(self.s_des[:, t].reshape(-1, 1) - current_state) + \
98             self.k_t
99         return action
```

Listing 4: Class for the P controller

```
102 class LQR_OPTIMAL(LQRSim):
103     def __init__(self) -> None:
104         super().__init__()
105         self.V_t = np.zeros_like(self.R_t)
106         self.V_t[-1] = self.R_t[-1]
107         self.v_t = np.zeros_like(self.r_t)
108         self.v_t[:, -1] = self.R_t[-1]@self.r_t[:, -1]
109         for t in range(self.T - 2, -1, -1):
110             M_t = (self.B_t/(self.H_t + self.B_t.T@self.V_t[t+1]@self.B_t))\
111                 @ self.B_t.T @ self.V_t[t + 1]@self.A_t
112             self.V_t[t] = self.R_t[t+1] + (self.A_t - M_t).T@self.V_t[t + 1]\
113                 @ self.A_t
114             self.v_t[:, t] = (
115                 self.R_t[t+1]@self.r_t[:, t+1].reshape(-1, 1)
116                 + (self.A_t - M_t).T
117                 @ (self.v_t[:, t + 1].reshape(-1, 1) - self.V_t[t + 1]
118                     @ self.b_t.reshape(-1, 1))
119                 ).reshape(-1)
120
121     def step_func(self, t):
122         current_state = self.s_hist[:, t].reshape(-1, 1)
123         action = -(self.H_t+self.B_t.T@self.V_t[t]@self.B_t)**-1\
124             * self.B_t.T\
125             @ (self.V_t[t]
126                 @ (self.A_t@current_state + self.b_t.reshape(-1, 1))
127                 - self.v_t[:, t].reshape(-1, 1))
128         return action
```

Listing 5: Class for the optimal controller