

# Optimal and Learning Control for Robotics

## Exercise Series 3

Bo Peng, New York University

Spring 2020

### Exercise 1

**Solution** a) The python code to run MPC at every time step is shown here:

```
1  def MPCstep(walking_model, xt, plan_horizon, Q_nominal, R_nominal, foot_position):
2      """
3      Replan at time step t with state xt.
4
5      Args:
6          walking_model: the environment;
7          xt: the initial state at time t;
8          plan_horizon: the replanning horizon length;
9          Q_nominal: nominal Q matrix;
10         R_nominal: nominal R matrix;
11         foot_position: a list of foot positions;
12
13     Returns:
14         ut: the control output at time step t
15     """
16
17     ## initialize matrices
18     G_bounds = []
19     h_bounds = []
20     Q = []
21     q = []
22     R = []
23     r = []
24
25     for i in range(plan_horizon):
26         Q.append(Q_nominal)
27         R.append(R_nominal)
28
29         q.append(Q_nominal.dot(np.array([-foot_position[i]], [0.])))
30         r.append(R_nominal.dot(np.array([-foot_position[i]])))
31
32         G_bounds.append(np.array([[0,0,1],[0,0,-1]]))
33         h_bounds.append(
34             np.array([[walking_model.foot_size+foot_position[i]],
35                       [walking_model.foot_size-foot_position[i]]]))
36
37     ## solve mpc with given codes
38     x_plan, u_plan = solve_mpc_collocation(
39         walking_model.A,
40         walking_model.B,
```

```

41         Q,
42         q,
43         R,
44         r,
45         G_bounds,
46         h_bounds,
47         plan_horizon,
48         xt)
49
50     ## take the first step
51     ut = u_plan[:, 0]
52
53     return ut
54
55
56 def MPCController(walking_model, horizon_length,
57                  plan_horizon, Q_nominal, R_nominal, foot_position):
58     """
59     A model predictive controller for stochastic walking environment.
60
61     Args:
62         walking_model: the environment;
63         horizon_length: the total length of horizon;
64         plan_horizon: the replanning horizon length;
65         Q_nominal: nominal Q matrix;
66         R_nominal: nominal R matrix;
67         foot_position: a list of foot positions;
68
69     Returns:
70         controller: a controller function for MPC.
71     """
72
73     def controller(x, i):
74         horizon = min(plan_horizon, horizon_length-i)
75         foot_steps = foot_position[i:i+horizon]
76
77         u = MPCstep(walking_model, x, horizon, Q_nominal, R_nominal, foot_steps)
78         return u
79
80     return controller
81

```

Listing 1: MPC Algorithm

The MPCController function returns a MPC controller that runs MPCstep function at each time step, which solve the original OC problem for a shorter horizon (plan\_horizon) and apply only the first control output  $u_t$ .

b) To find the shortest time horizon I run an experiment that use different plan\_horizon to solve the OC problem. In each experiment, the same plan\_horizon is used for 200 times since we are using a noisy environment.

To check if the controller stabilizes the system, I define a metric "cost", which is the distance between original plan states and real states:

$$\text{cost} = \sqrt{(x_{\text{real}} - x_{\text{plan}})^T (x_{\text{real}} - x_{\text{plan}})} \quad (1)$$

The idea is that if the controller could stabilize the system in a noisy environment, then the system should behave as if it is in a non-noise environment, which is the planned state. Therefore, if the cost is low, then the controller stabilize the system in the noisy environment.

The experiment results are shown below (note that the experiment results may change because of the noise in the environment):

Table 1: Plan Horizon Experiment Results

Plan Horizon	Mean Cost	Minimum Cost	Maximum Cost
1	2.059134e+09	2.154437e+09	1.953782e+09
2	1.965600e+09	2.074292e+09	1.872367e+09
3	1.865734e+09	1.953721e+09	1.780687e+09
4	1.714734e+09	1.778597e+09	1.630459e+09
5	1.509684e+09	1.581568e+09	1.420585e+09
6	1.241064e+09	1.288699e+09	1.167756e+09
7	9.014475e+08	9.564713e+08	8.431085e+08
8	5.677452e+08	6.328668e+08	5.134207e+08
9	3.008015e+08	3.526146e+08	2.464213e+08
10	9.021696e+07	1.335314e+08	3.131681e+07
11	1.306313e+05	4.278392e+06	9.459968e-01
12	1.954138e+05	1.954131e+07	6.288541e-01
13	4.634063e+05	2.165111e+07	4.381214e-01
14	3.763186e-01	6.592180e-01	3.154441e-01
15	2.979701e-01	4.668112e-01	2.368805e-01
16	2.418781e+05	1.303302e+07	1.844776e-01
17	1.759316e+04	1.759294e+06	1.379911e-01
18	6.661043e+04	6.661023e+06	1.092199e-01
19	1.851220e-01	5.788336e-01	9.333695e-02
20	1.702376e-01	3.820062e-01	9.580803e-02

If a minimum cost is closer to 0, it means that the controller could stabilize the system in at least 1 trails for the used plan horizon. So the shortest time horizon is 11. Note that it is very likely to fail if we use this plan horizon directly. A failure example is shown below, and the animation can be found [here](#) (the planned curve is the original planned one in the non-noise environment):

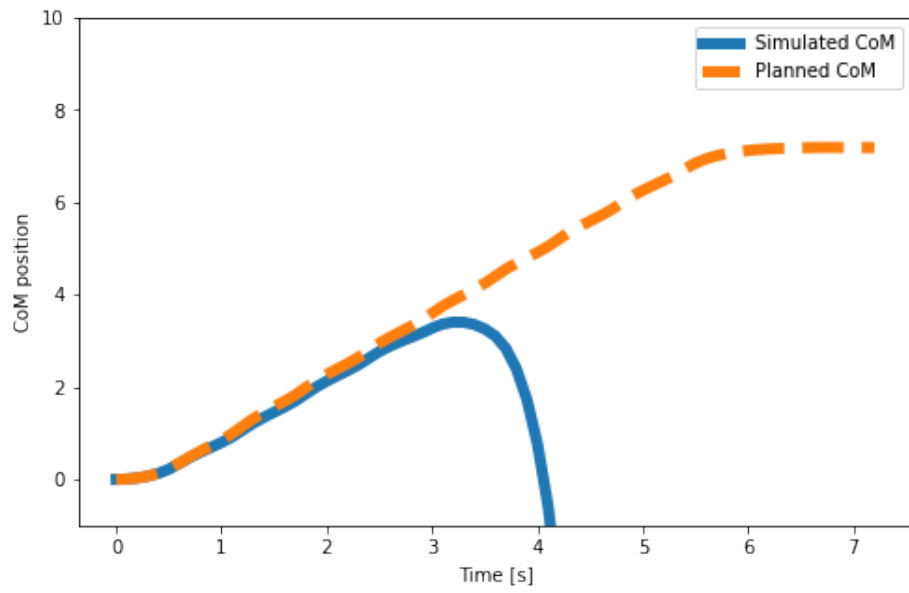


Figure 1: CoM Position (horizon=11)

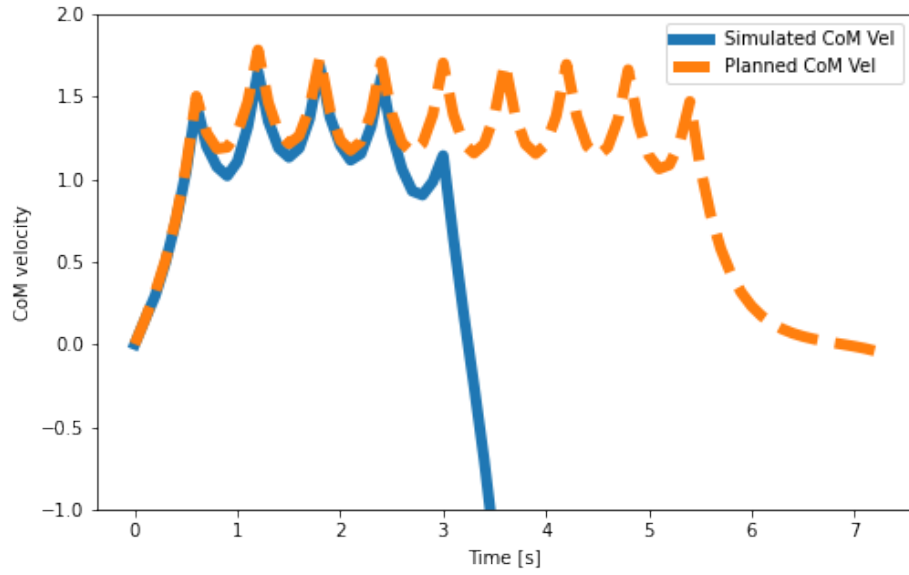


Figure 2: CoM Velocity (horizon=11)

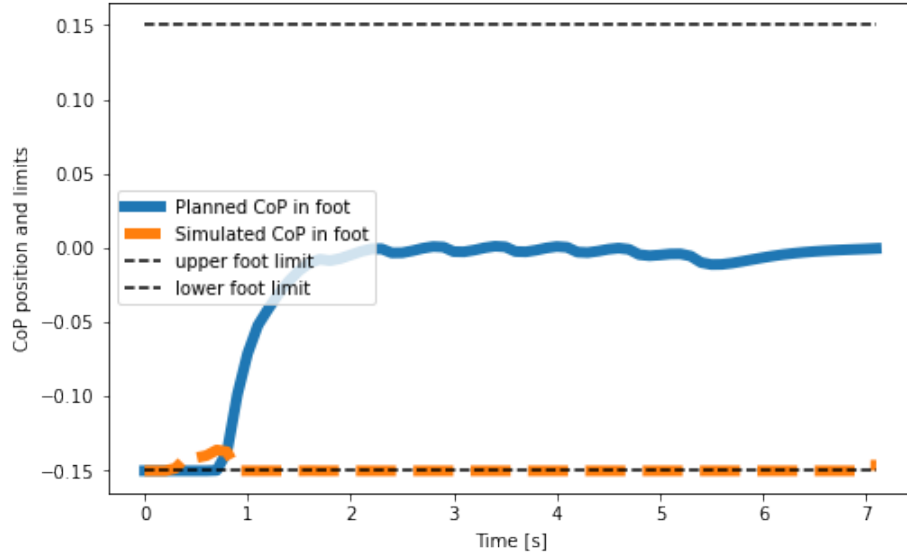


Figure 3: CoP Position and Limits (horizon=11)

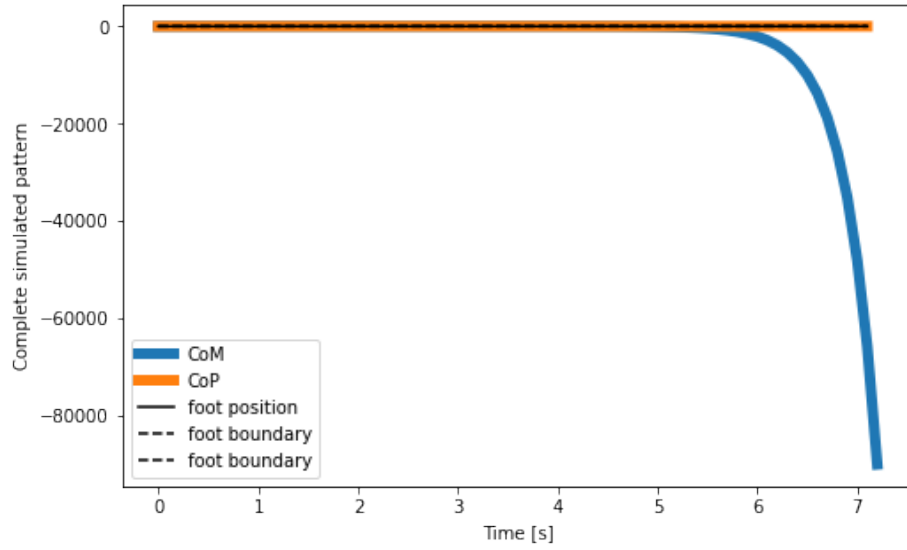


Figure 4: Complete Simulated Pattern (horizon=11)

A successful example is shown below, and the animation can be found [here](#).

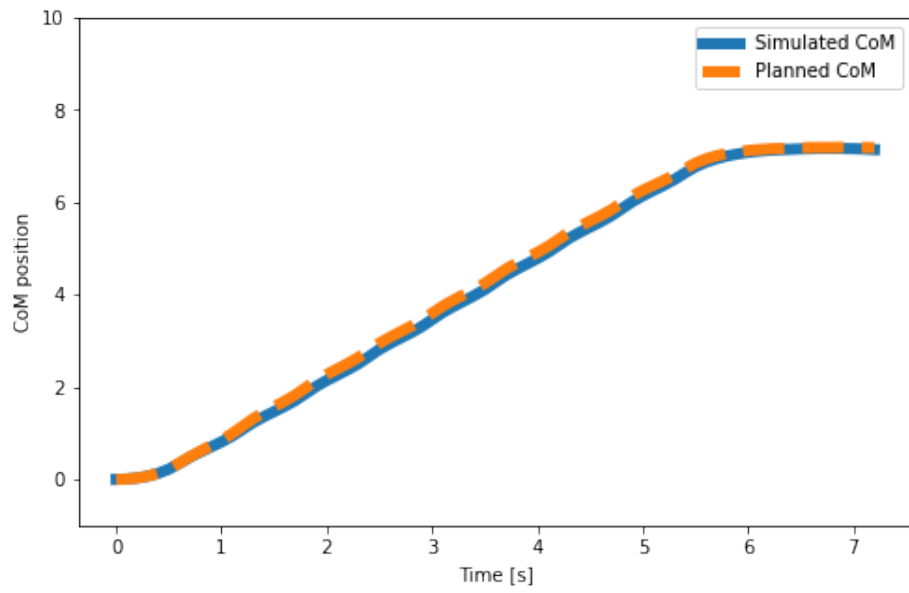


Figure 5: CoM Position (horizon=11)

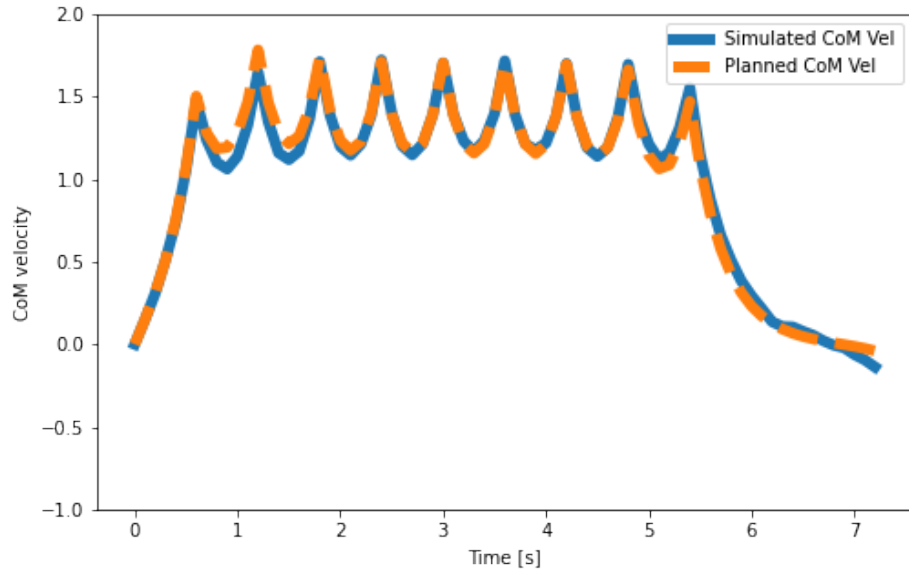


Figure 6: CoM Velocity (horizon=11)

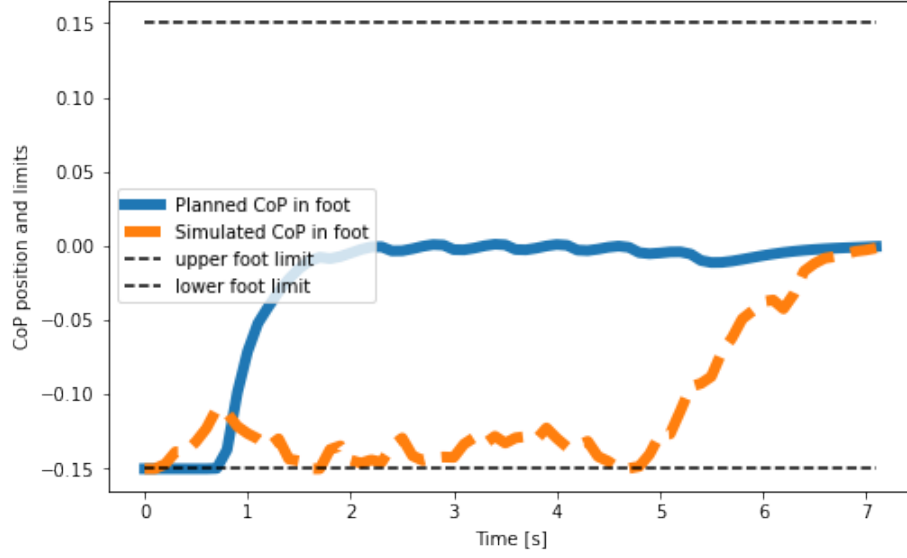


Figure 7: CoP Position and Limits (horizon=11)

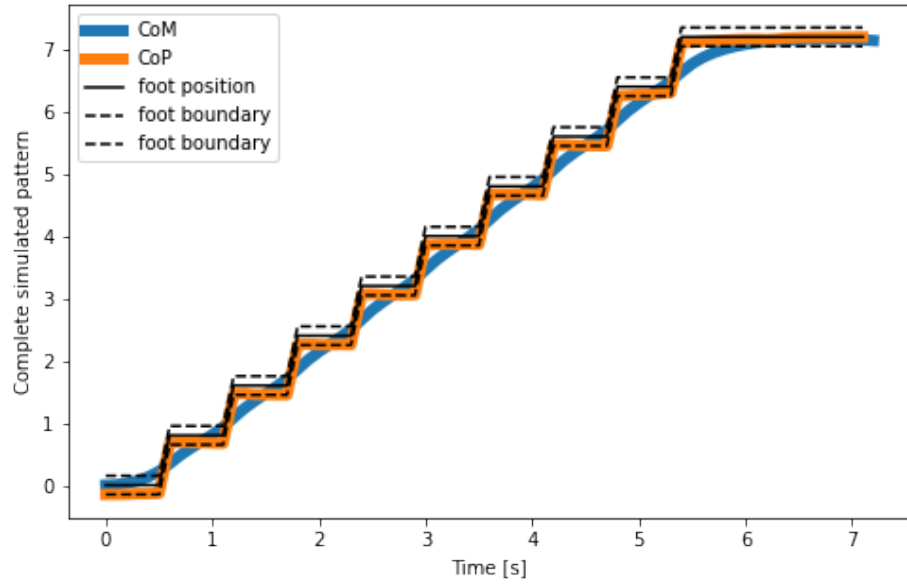


Figure 8: Complete Simulated Pattern (horizon=11)

c) If a mean cost is closer to 0, it means that the controller could stabilize the system in most of the trails. So a stable controller should plan at least 14 steps in advance. An example is shown below, and the

animation can be found [here](#).

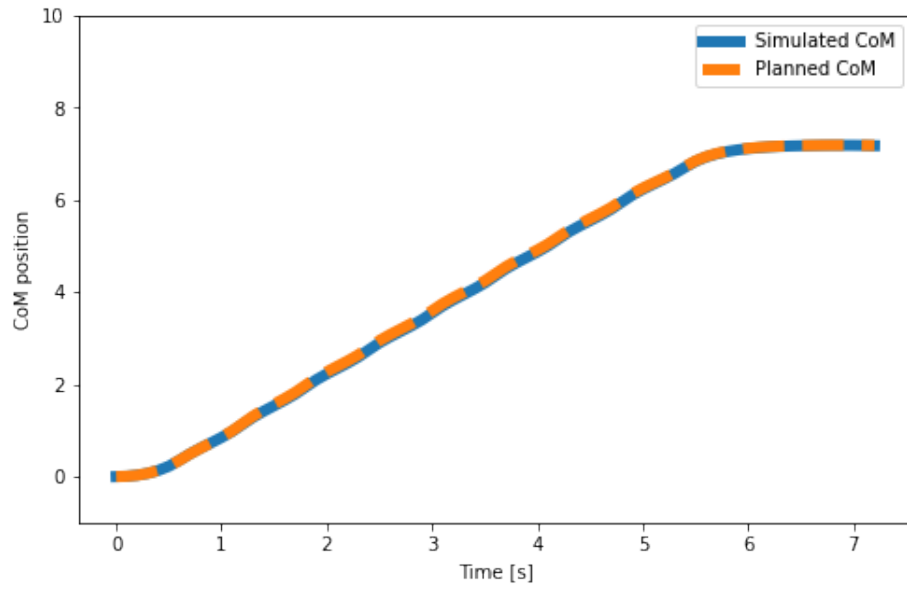


Figure 9: CoM Position (horizon=14)

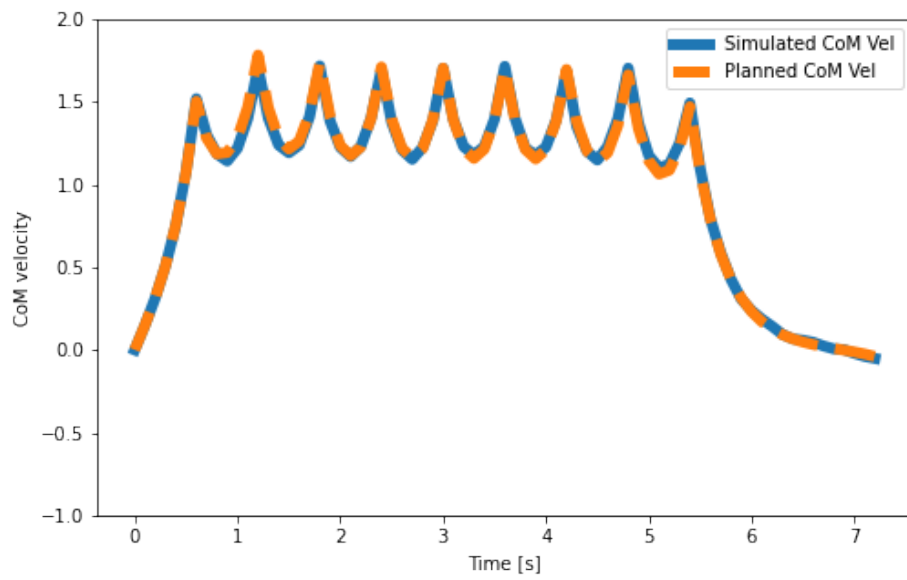


Figure 10: CoM Velocity (horizon=14)



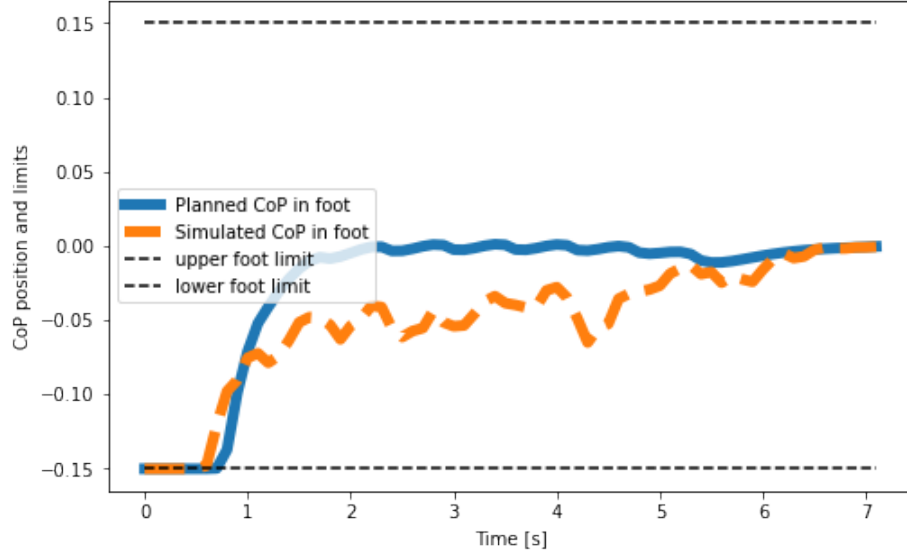


Figure 11: CoP Position and Limits (horizon=14)

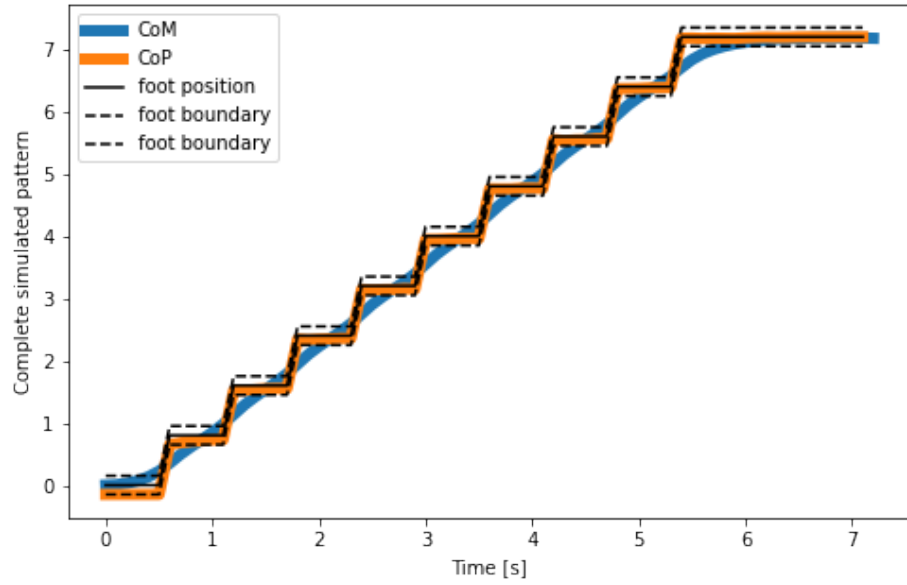


Figure 12: Complete Simulated Pattern (horizon=14)

d) The comparison between the original planned control and executed control is shown below. In general, the executed control is a little bit lower than the planned control sequence to adjust to the environment

noise. The longer plan horizon, the smaller difference there will be.

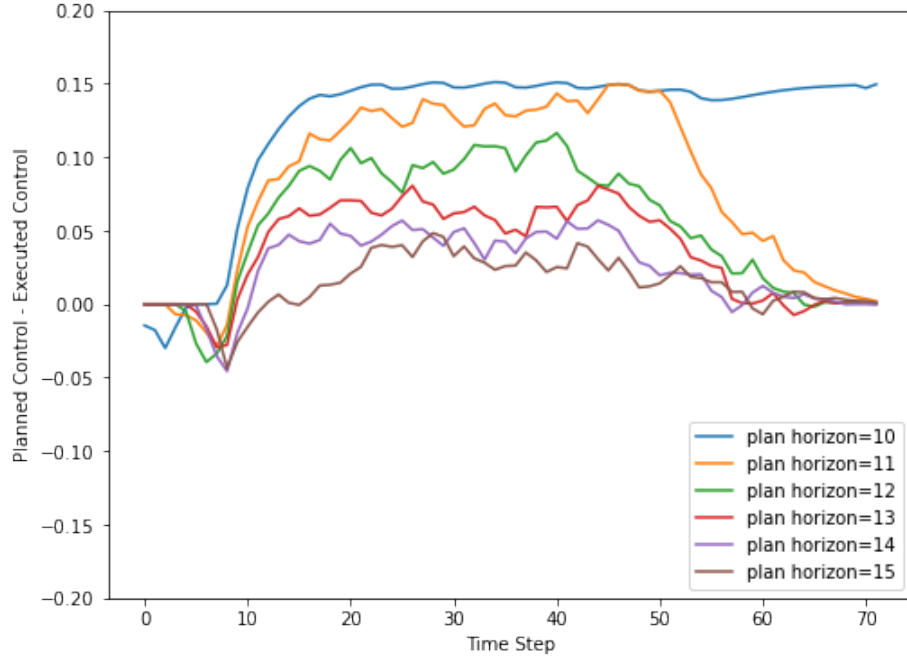


Figure 13: Comparison between Control Sequence

e) When I increase the cost weight in cost function, there's no big difference in shortest plan horizon or stable plan horizon. Besides, if I multiply the cost weight by 1,000, the controller fails to guide the robot to walk. I think maybe that's because when we increase the cost weight the controller will force the robot to follow the trajectory, which maybe very hard in the noisy environment. If the cost weight is too high, it could be impossible to follow the trajectory and lead to failure.

Similarly, if decrease the cost weight there's no big difference either. If I divide the cost weight by 1,000, the controller also fails to guide the robot. I think that's because when we decrease the weight, the constraints on the target will become less important. So the robot won't need to follow the trajectory to get low cost. In conclusion, only change the cost weight may not decrease the plan horizon, and may lead to failure in extreme cases.

When I increase the terminal cost, the shortest control horizon decrease to 7. The simulation results are shown below, and the animation can be found [here](#). It seems that if we increase terminal cost, the controller would force the robot to stay at the terminal, which is helpful for planning.

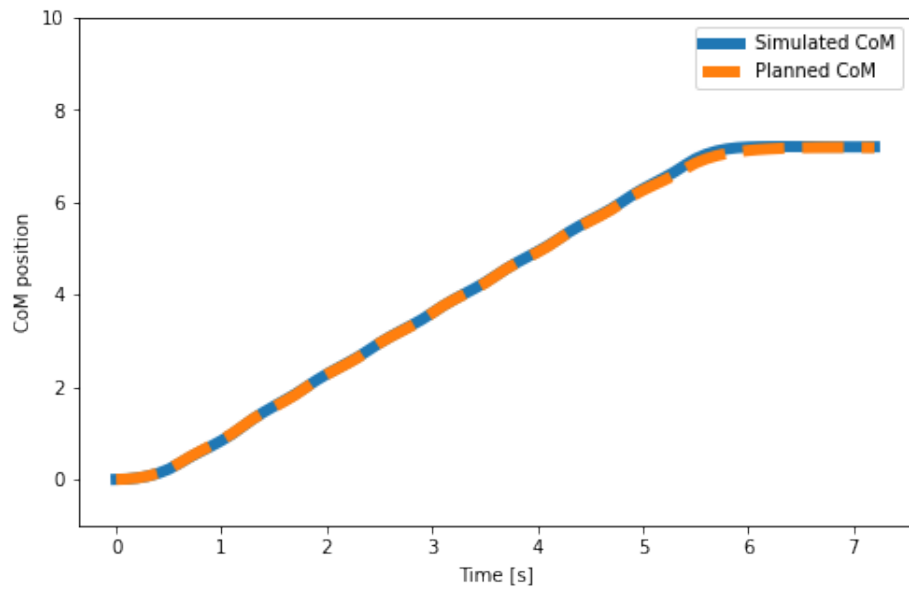


Figure 14: CoM Position (horizon=7)

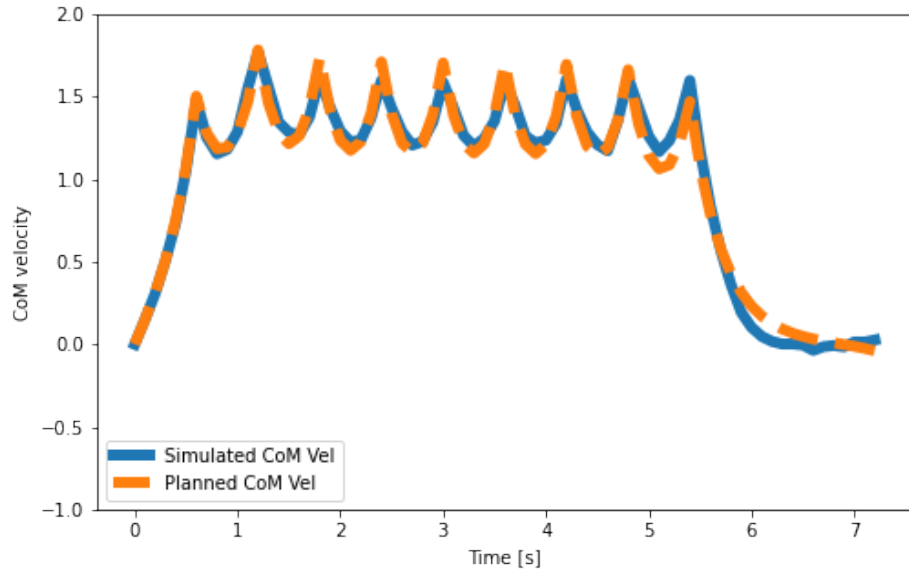


Figure 15: CoM Velocity (horizon=7)

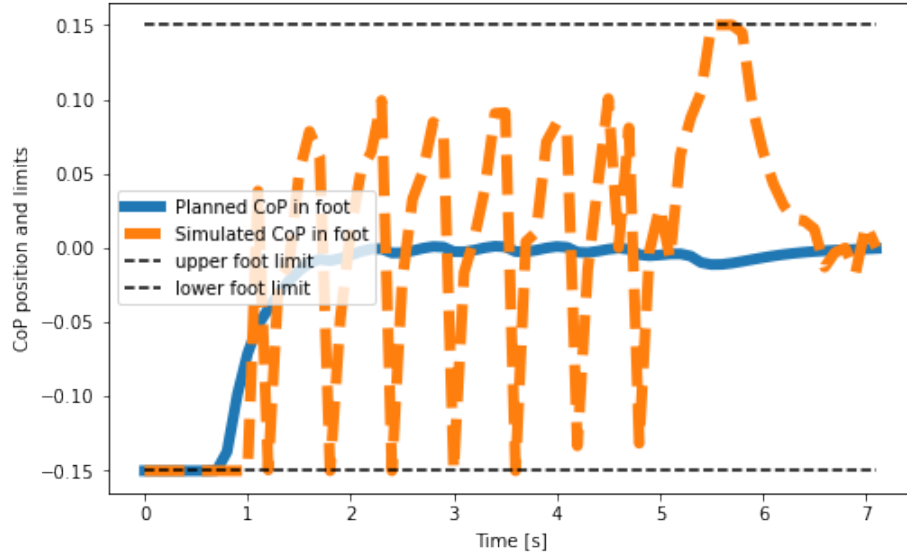


Figure 16: CoP Position and Limits (horizon=7)

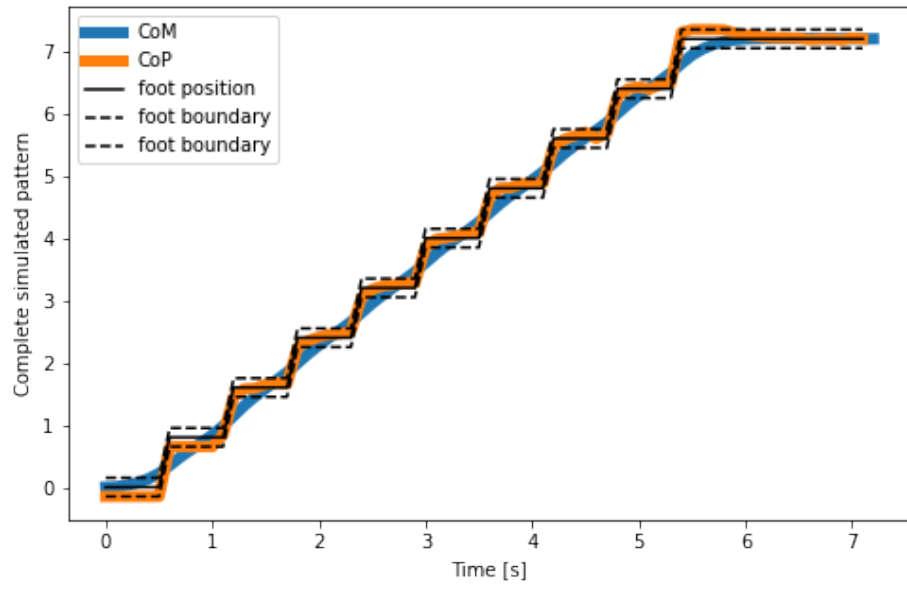


Figure 17: Complete Simulated Pattern (horizon=7)

## Exercise 2

**Solution** a) To run the value iteration algorithm, first we need to code the environment. The environment is shown below:

```
1  def env(x, u):
2      """
3      The grid-world environemnt.
4
5      Args:
6          x: the current state;
7          u: the current action.
8
9      Returns:
10         x_next: the next state;
11         cost: the current cost.
12     """
13
14     def get_coord(x):
15         """
16         A helper function to transfer state x to coordinate in grid-world.
17
18         Args:
19             x: the current.
20
21         Returns:
22             i, j: the coordinate in grid-world.
23         """
24
25         return x//5, x%5
26
27     i, j = get_coord(x)
28     cost = 0
29
30     if (i==0 and j==4) or (i==3 and j==4):
31         cost = -1
32     elif (i==0 and j==2) or (i==1 and j==2) or (i==2 and j==3) or (i==3 and j==0):
33         cost = 1
34     elif (i==1 and j==4) or (i==2 and j==1) or (i==3 and j==3):
35         cost = 10
36
37     di, dj = 0, 0
38     if u == 0:
39         di = -1
40     elif u == 1:
41         di = 1
42     elif u == 2:
43         dj = -1
44     elif u == 3:
45         dj = 1
46
47     ## obstacles
48     if (i+di==1 and j+dj==1) or (i+di==2 and j+dj==2):
49         di, dj = 0, 0
50
51     ## make a move
52     i += di
53     j += dj
54
55     ## boundry condition
```

```

56     if i < 0:
57         i = 0
58     elif i > 3:
59         i = 3
60
61     if j < 0:
62         j = 0
63     elif j > 4:
64         j = 4
65
66     x_next = int(i * 5 + j)
67
68     return x_next, cost
69

```

Listing 2: Environment Dynamics

The value iteration algorithm is shown below:

```

1  def ValueIteration(env, states, actions, alpha, tol=1e-6):
2      """
3      Value iteration algorithm.
4
5      Args:
6          env: the environment function;
7          states: the state indices;
8          actions: the action indices;
9          alpha: the discount factor;
10         tol: converge tolerance.
11
12     Returns:
13         v: the value function for each state.
14     """
15
16     v = np.zeros_like(states, dtype=np.float)
17     i = 0
18     while True:
19         v_new = np.zeros_like(v)
20         for s in states:
21             q = np.zeros_like(actions, dtype=np.float)
22             for a in actions:
23                 s_next, cost = env(s, a)
24                 q[a] = cost + alpha * v[s_next]
25
26             v_new[s] = np.min(q)
27
28         if np.all(np.abs(v_new-v) < tol):
29             break
30         v = v_new
31         i += 1
32
33     print(f"Converge in {i} iterations.")
34
35     return v
36

```

Listing 3: Value Iteration Algorithm

In this grid-world case, the algorithm converge in 1,375 iterations. The optimal value function and optimal policy are shown below:

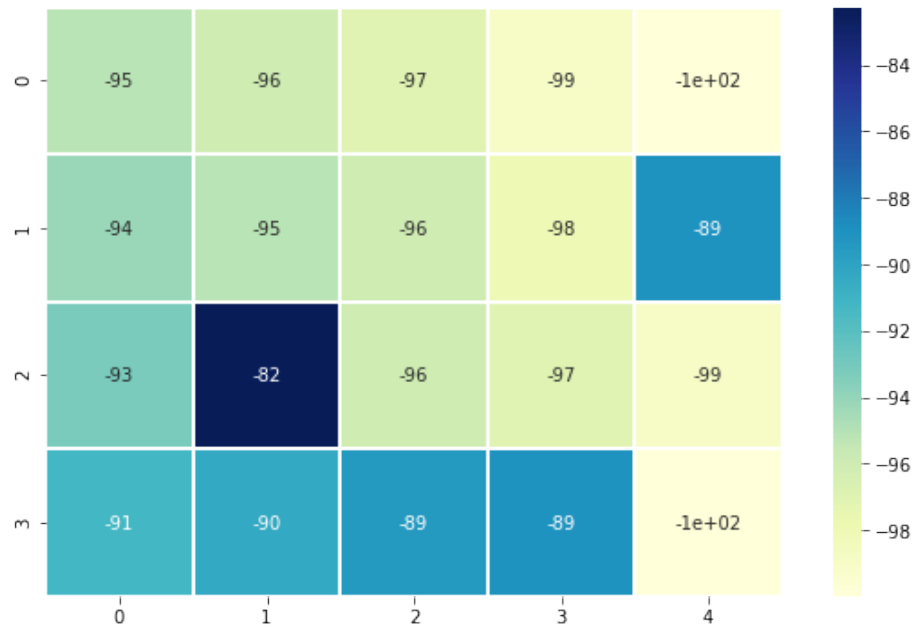


Figure 18: Optimal Value Function (Value Iteration)

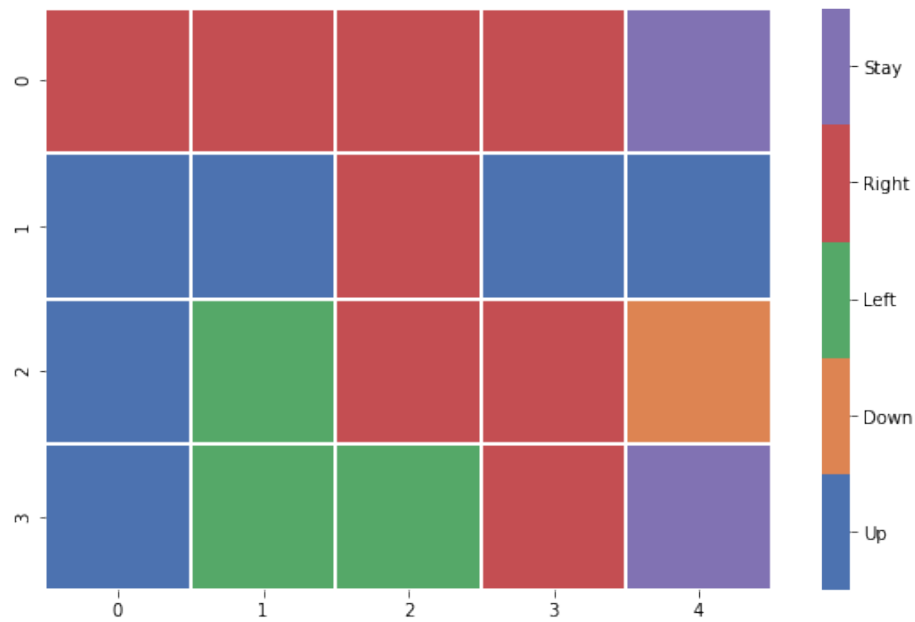


Figure 19: Optimal Policy (Value Iteration)

2) The policy iteration algorithm is shown below:

```
1  def PolicyIteration(env, states, actions, policy, alpha):
2      """
3      Policy iteration algorithm.
4
5      Args:
6          env: the environment function;
7          states: the state indices;
8          actions: the action indices;
9          policy: the initial policy;
10         alpha: the discount factor;
11         tol: converge tolerance.
12
13     Returns:
14         v: the value function for each state.
15     """
16
17     def PolicyEvaluate(policy):
18         """
19         A helper function to evaluate the given policy.
20         """
21
22         N = len(states)
23         T = np.zeros((N, N))
24         g = np.zeros_like(states)
25
26         for s in states:
27             a = policy[s]
28             s_next, g[s] = env(s, a)
29             T[s, s_next] = 1
30
31         J = np.linalg.inv(np.eye(N) - alpha * T) @ g
32
33         return J
34
35     def PolicyUpdate(policy, J):
36         """
37         A helper function to update current policy.
38         """
39
40         new_policy = np.zeros_like(policy)
41         for s in states:
42             q = np.zeros_like(actions, dtype=np.float)
43             for a in actions:
44                 s_next, cost = env(s, a)
45                 q[a] = cost + alpha * J[s_next]
46
47             new_policy[s] = np.argmin(q)
48
49         return new_policy
50
51     i = 0
52     while True:
53         ## policy evaluation
54         J = PolicyEvaluate(policy)
55         ## policy update
56         new_policy = PolicyUpdate(policy, J)
57
58         if np.all(policy == new_policy):
59             return J
```



```

60         break
61
62         policy = new_policy
63         i += 1
64
65     print(f"Converge in {i} iterations.")
66
67     return policy, J
68

```

Listing 4: Policy Iteration

In this grid-world case, the algorithm converge in 8 iterations. The optimal value function and optimal policy are shown below:

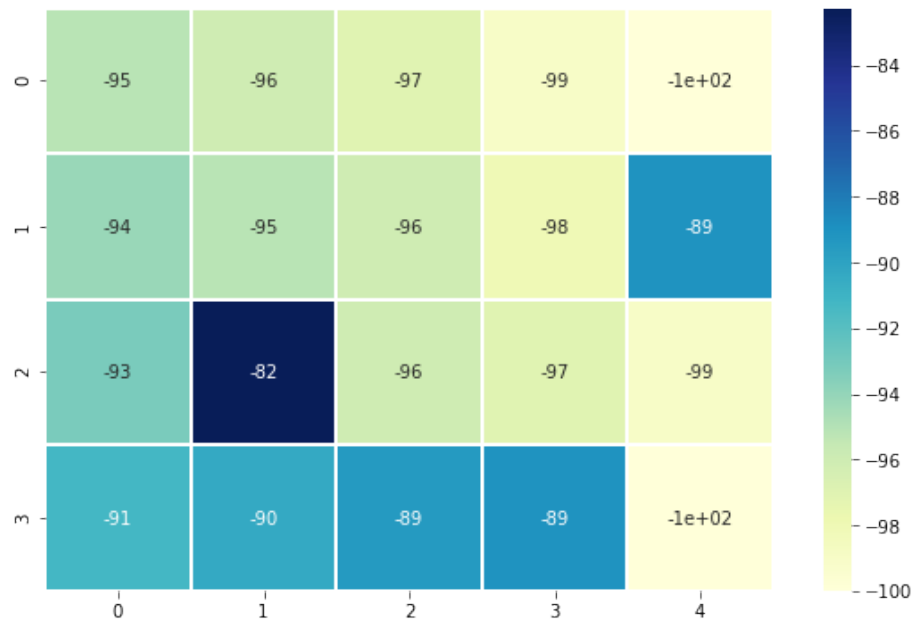


Figure 20: Optimal Value Function (Policy Iteration)

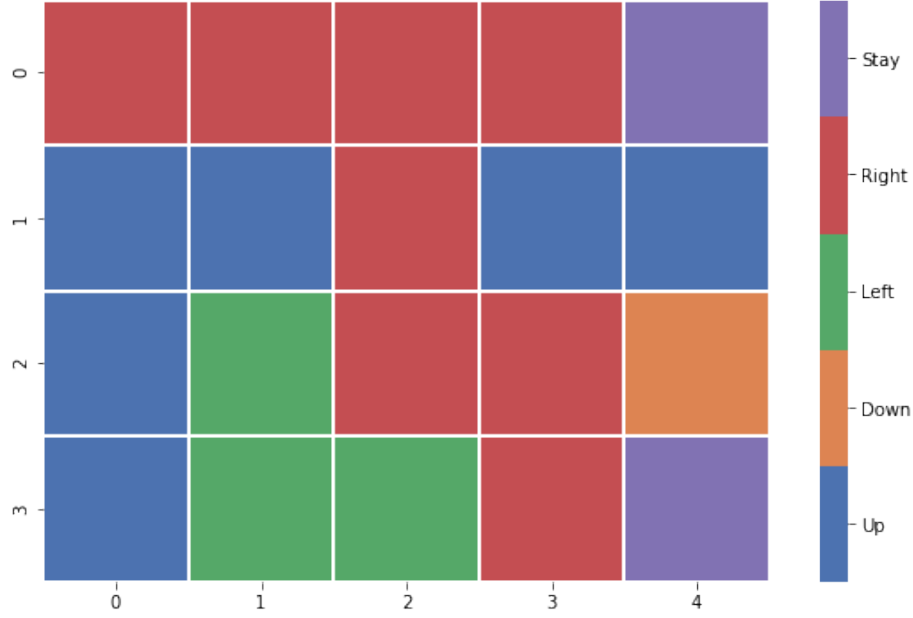


Figure 21: Optimal Policy (Policy Iteration)

c) For the solution, both value iteration and policy iteration algorithm could find the optimal value function/ policy. The optimal value function/ policy are the same in this case, and they are already shown in the previous sections.

For the convergence/ complexity, both value iteration and policy iteration algorithm could converge, and the corresponding optimal value/ policy are the same. But the convergence speed is different. The value iteration algorithm may take a relative longer time to converge, in this case is 1,375, while policy iteration only needs 8 iterations. So policy iteration could accelerate the convergence. Note that we are using the linear equation when we evaluate the policy. If we use the iterative method, the policy iteration algorithm may need longer time to converge.