

The Solution of Model Predictive Control: Theory,
Computation, and Design [1]

lixc21

January 1, 2023

Contents

1	Getting Started with Model Predictive Control	2
1.1	Brief Review	2
1.2	The Solution of Exercises	2

Chapter 1

Getting Started with Model Predictive Control

1.1 Brief Review

In this section, we just consider state space linear time invariant system with zero steady state.

Lemma 1.3 (LQR convergence). For (A, B) controllable, the infinite LQR gives a convergent closed-loop system.

Proof. Because (A, B) is controllable, there exists a sequence of n inputs that transfers the state to zero. When $k > n$, we let $u = 0$, then the objective function $V(x, u) = \sum_{k=0}^{\infty} x_k^T Q x_k + u^T R u$ is finite, which implies the optimization problem is feasible. On the other hand, the solution is unique since $R > 0$ and the objective function is strict convex with u .

So the solution of the LQR problem exists and is unique. This implies to that the objective function is non-increasing with time, and we have $x \rightarrow 0, u \rightarrow 0$ as $k \rightarrow \infty$. \square

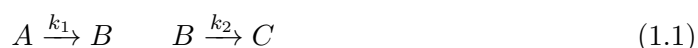
Remark. The optimal solution can be calculate from Riccati equation, which is from backward dynamic programming similar to Kalman filter.

$$\begin{aligned} K &= -(B^T P B + R)^{-1} B^T P A \\ P &= Q + A^T P A - A^T P B (B^T P B + R)^{-1} B^T P A \end{aligned}$$

1.2 The Solution of Exercises

Exercise 1.1. State space form for chemical reaction model.

Consider the following chemical reaction kinetics for a two-step series reaction



We wish to follow the reaction in a constant volume, well-mixed, batch reactor. As taught in the undergraduate chemical engineering curriculum, we proceed by writing material balances for the three species giving

$$\frac{dc_A}{dt} = -r_1 \quad \frac{dc_B}{dt} = r_1 - r_2 \quad \frac{dc_C}{dt} = r_2 \quad (1.2)$$

in which c_j is the concentration of species j , and r_1 and r_2 are the rates (mol/(time·vol)) at which the two reactions occur. We then assume some rate law for the reaction kinetics, such as

$$r_1 = k_1 c_A \quad r_2 = k_2 c_B \quad (1.3)$$

We substitute the rate laws into the material balances and specify the starting concentrations to produce three differential equations for the three species concentrations.

- (a) write the linear state space model for the deterministic series chemical reaction model. Assume we can measure the component A concentration. What are x , y , A , B , C , and D for this model?
- (b) Simulate this model with initial conditions and parameters given by

$$c_{A0} = 1 \quad c_{B0} = c_{C0} = 0 \quad k_1 = 2 \quad k_2 = 1$$

Answer 1. (a) the linear state space model is

$$\frac{dx}{dt} = \begin{bmatrix} -k_1 & 0 & 0 \\ k_1 & -k_2 & 0 \\ 0 & k_2 & 0 \end{bmatrix} x = Ax \quad (1.4)$$

where $x = [c_A, c_B, c_C]^\top$. B does not exist because there is no system input variables. $C = [1, 0, 0]^\top$, $D = 0$, $y = Cx$.

- (b) the simulation result is shown as Fig.1.1. The code used in all of the exercise can be found in github <https://github.com/lixc21/MPC-Solution>.

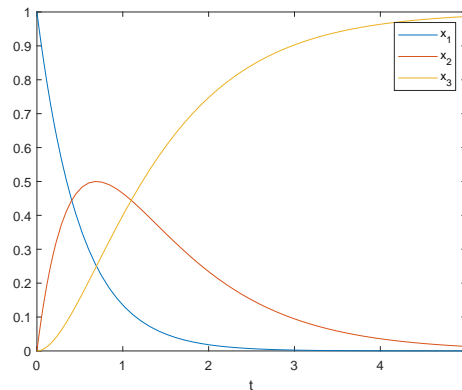


Figure 1.1: system simulation

Exercise 1.2. Distributed systems and time delay.

We assume familiarity with the transfer function of a time delay from an undergraduate systems course

$$\bar{y}(s) = e^{-\theta s} \bar{u}(s) \quad (1.5)$$

Let's see the connection between the delay and the distributed systems, which give rise to it. A simple physical example of a time delay caused by transport in a flowing system. Consider plug flow in a tube depicted in Fig.1.2.

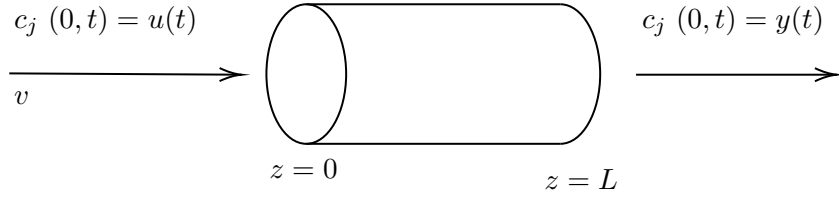


Figure 1.2: Plug-flow reactor

- (a) Write down the equation of change for moles of component j for an arbitrary volume element and show that

$$\frac{\partial c_j}{\partial t} = -\nabla \cdot (c_j v_j) + R_j \quad (1.6)$$

in which c_j is the molar concentration of component j , v_j is the velocity of component j , and R_j is the production rate of component j due to chemical reaction. Plug flow means the fluid velocity of all components is purely in the z direction, and is independent of r and θ and, we assume here, z

$$v_j = v \delta_z \quad (1.7)$$

- (b) Assuming plug flow and neglecting chemical reaction in the tube, show that the equation of change reduces to

$$\frac{\partial c_j}{\partial t} = -v \frac{\partial c_j}{\partial z} \quad (1.8)$$

This equation is known as a hyperbolic, first-order partial differential equation.

$$c_j(z, t) = u(t) \quad 0 = z \quad t \geq 0 \quad (1.9)$$

$$c_j(z, t) = c_{j0}(t) \quad 0 \leq z \leq L \quad t = 0 \quad (1.10)$$

In other words, we are using the feed concentration as the manipulated variable, $u(t)$, and the tube starts out with some initial concentration profile of component j , $c_{j0}(z)$.

- (c) Show that the solution to (1.8) with these boundary conditions is

$$c_j(z, t) = \begin{cases} u(t - z/v) & vt > z \\ c_{j0}(z - vt) & vt < z \end{cases} \quad (1.11)$$

- (d) If the reactor starts out empty of component j , show that the transfer function between the outlet concentration, $y = c_j(L, t)$, and the inlet concentration, $c_j(0, t) = u(t)$, is a time delay. What is the value of θ ?

Answer 2. (a) let f be the moles of one of the component, then from 3D Leibniz formula, we get

$$\frac{\partial c_j}{\partial t} = \frac{d}{dt} \int_V f(\vec{x}, t) dV = \int_V \frac{\partial f}{\partial t} dV - \int_A f \vec{v} \cdot \vec{n} dV \quad (1.12)$$

where V is a small unit volume, A is the responding surface, v represent the velocity of the point on the surface, n is the outward unit normal vector related

to u .

By using Gauss divergence theorem, we know that

$$\int_V \frac{\partial f}{\partial t} dV - \int_A f \vec{v} \cdot \vec{n} dV = \int_V \frac{\partial f}{\partial t} - \nabla \cdot f \vec{v} dV = -\nabla \cdot (c_j v_j) + R_j \quad (1.13)$$

where the last equation comes from $v_j = v \delta_z$.

- (b) Neglecting chemical reaction in the tube, we get $R_j = 0$. Then we know that

$$\frac{\partial c_j}{\partial t} = -\nabla \cdot (c_j v_j) = -\left(\frac{\partial}{\partial x} \delta_x + \frac{\partial}{\partial y} \delta_y + \frac{\partial}{\partial z} \delta_z \right) \cdot v \delta_z = -v \frac{\partial c_j}{\partial z} \quad (1.14)$$

- (c) Assuming that $u(t - z/v) = c_{j0}(z - vt)$ when $vt < z$, we just need to prove the solution is $c_j(z, t) = u(t - z/v)$. The variables of original partial differential equation has already been separated, so we get $c_j(z, t) = u(t - z/v)$ easily from the method of characteristics.

- (d) We know that $y = u(t - L/v)$, which is a time delay. The value of θ could be L/v .

Exercise 1.3. Pendulum in the state space.

Consider the pendulum suspended at the end of a rigid link depicted in Figure 1.3. Let r and θ denote the polar coordinates of the center of the pendulum, and let $p = r \delta_r$ be the position vector of the pendulum, in which δ_r and δ_θ are the unit vectors in polar coordinates. We wish to determine a state space description of the system. We are able to apply a torque T to the pendulum as our manipulated variable. The pendulum has mass m , the only other external force acting on the pendulum is gravity, and we neglect friction. The link provides force $-t \delta_r$ necessary to maintain the pendulum at distance $r = R$ from the axis of rotation, and we measure the force t .

- (a) Provide expressions for the four partial derivatives for changes in the unit vectors with r and θ

$$\frac{\partial \delta_r}{\partial r} \quad \frac{\partial \delta_r}{\partial \theta} \quad \frac{\partial \delta_\theta}{\partial r} \quad \frac{\partial \delta_\theta}{\partial \theta} \quad (1.15)$$

- (b) Use the chain rule to find the velocity of the pendulum in terms of the time derivatives of r and θ . Do not simplify yet by assuming r is constant. We want the general result.

- (c) Differentiate again to show that the acceleration of the pendulum is

$$\ddot{p} = (\ddot{r} - r\dot{\theta}^2)\delta_r + (r\ddot{\theta} + 2\dot{r}\dot{\theta})\delta_\theta \quad (1.16)$$

- (d) Use a momentum balance on the pendulum mass (you may assume it is a point mass) to determine both the force exerted by the link

$$t = mR\dot{\theta}^2 + mg \cos \theta \quad (1.17)$$

and an equation for the pendulum due to gravity and the applied torque

$$mR\ddot{\theta} - T/R + mg \sin \theta = 0 \quad (1.18)$$

- (e) Define a state vector and give a state space description of your system. What is the physical significance of your state. Assume you measure the force exerted by the link.

One answer is

$$\frac{dx_1}{dt} = x_2 \quad (1.19)$$

$$\frac{dx_2}{dt} = -(g/R) \sin x_1 + u \quad (1.20)$$

$$y = mRx_2^2 + mg \cos x_1 \quad (1.21)$$

in which $u = T/(mR)$

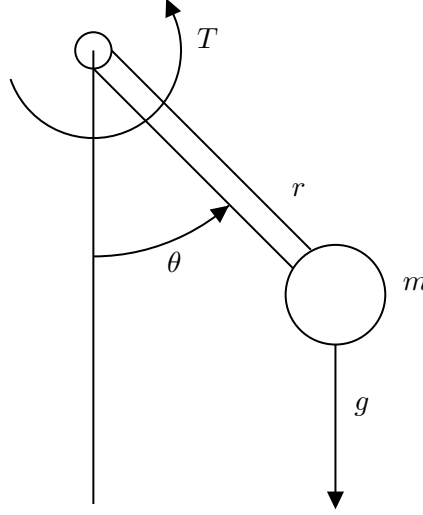


Figure 1.3: Pendulum with applied torque

Answer 3. (a) we assume that δ_θ is rotated from δ by anticlockwise.

$$\frac{\partial \delta_r}{\partial r} = 0 \quad \frac{\partial \delta_r}{\partial \theta} = \delta_\theta \quad \frac{\partial \delta_\theta}{\partial r} = 0 \quad \frac{\partial \delta_\theta}{\partial \theta} = -\delta_r \quad (1.22)$$

(b) Since $p = r\delta_r$

$$\dot{p} = \dot{r}\delta_r + r\frac{\partial \delta_r}{\partial t} = \dot{r}\delta_r + r\frac{\partial \delta_r}{\partial \theta}\dot{\theta} = \dot{r}\delta_r + r\dot{\theta}\delta_\theta \quad (1.23)$$

(c) Differentiate again

$$\ddot{p} = \ddot{r}\delta_r + \dot{r}\dot{\theta}\delta_\theta + r\ddot{\theta}\delta_\theta + r\dot{\theta}(-\delta_r\dot{\theta}) + \dot{r}\dot{\theta}\delta_\theta \quad (1.24)$$

$$= (\ddot{r} - r\dot{\theta}^2)\delta_r + (r\ddot{\theta} + 2\dot{r}\dot{\theta})\delta_\theta \quad (1.25)$$

(d) By the Newton's second law of motion, we get

$$F = -t\delta_r + T/R\delta_\theta + mg \sin \theta \delta_r + mg \cos \theta \delta_\theta = -mR\dot{\theta}^2\delta_r \quad (1.26)$$

Simplify it by two direction

$$t = mR\dot{\theta}^2 + mg \cos \theta \quad (1.27)$$

$$mR\ddot{\theta} - T/R + mg \sin \theta = 0 \quad (1.28)$$

(e) State vector could be $x = [\theta, \dot{\theta}]^\top$, and the system

$$\dot{x} = \begin{bmatrix} \dot{\theta} \\ \ddot{\theta} \end{bmatrix} = \begin{bmatrix} \dot{\theta} \\ -g \sin \theta / R + T / (mR^2) \end{bmatrix} \quad (1.29)$$

$$y = \theta \quad (1.30)$$

Exercise 1.4. Time to Laplace domain.

Take the Laplace transform of the following set of differential equations and find the transfer function, $G(s)$, connecting $\bar{u}(s)$ and $\bar{y}(s)$, $\bar{y} = G\bar{u}$

$$\frac{dx}{dt} = Ax + Bu \quad (1.31)$$

$$y = Cx + Du \quad (1.32)$$

For $x \in \mathbb{R}^n$, $y \in \mathbb{R}^p$, and $u \in \mathbb{R}^m$, what is the dimension of the G matrix? What happens to the initial condition, $x(0) = x_0$?

Answer 4. The Laplace transform of the differential equation is

$$sx = Ax + Bu \quad (1.33)$$

and the transfer function

$$G(s) = C(sI - A)^{-1}B \in \mathbb{R}^{p \times m} \quad (1.34)$$

the initial condition does not appear in the Laplace transform, because the Laplace transform explains the dynamic from u to y , and when we need to determine the accurate trajectory of the system, the initial condition is needed by inverse Laplace transform.

Exercise 1.5. Converting between continuous and discrete time models.

Given a prescribed $u(t)$, derive and check the solution to (1.31). Given a prescribed $u(k)$ sequence, what is the solution to the discrete time model

$$x(k+1) = \tilde{A}x(k) + \tilde{B}u(k) \quad (1.35)$$

$$y(k) = \tilde{C}x(k) + \tilde{D}u(k) \quad (1.36)$$

- (a) Compute \tilde{A} , \tilde{B} , \tilde{C} , and \tilde{D} so that the two solutions agree at the sample times for a zero-order hold input, i.e., $y(k) = y(t_k)$ for $u(t) = u(k)$, $t \in (t_k, t_{k+1})$ in which $t_k = k\Delta$ for sample time Δ .
- (b) Is your result valid for A singular? If not, how can you find \tilde{A} , \tilde{B} , \tilde{C} , and \tilde{D} for this case?

Answer 5. (a) the solution to (1.31) is

$$x(\Delta) = e^{A\Delta}x_0 + \int_0^\Delta e^{A(\Delta-\tau)}Bu(\tau)d\tau \quad (1.37)$$

so the accurate discrete time model is

$$\tilde{A} = e^{A\Delta} \quad \tilde{B} = \int_0^\Delta e^{A(\Delta-\tau)}Bd\tau \quad \tilde{C} = C \quad \tilde{D} = D \quad (1.38)$$

- (b) Yes, this result can be calculate normally even if A is singular.

Exercise 1.6. Continuous to discrete time conversion for nonlinear models

Consider the autonomous nonlinear differential equation model

$$\frac{dx}{dt} = f(x, u) \quad (1.39)$$

$$x(0) = x_0 \quad (1.40)$$

Given a zero-order hold on the input, let $s(t, u, x_0)$, $0 \leq t \leq \Delta$, be the solution to (1.39) given initial condition x_0 at the time $t = 0$, and constant input u is applied for t in the interval $0 \leq t \leq \Delta$. Consider also the nonlinear discrete time model

$$x(k+1) = F(x(k), u(k)) \quad (1.41)$$

- (a) What is the relationship between F and s so that the solution of the discrete time model agrees at the sample times with the continuous time model with a zero-order hold?
- (b) Assume f is linear and apply this result to check the result of Exercise 1.5.

Answer 6. (a) The relationship is

$$F(x(k), u(k)) = s(\Delta, u(k-1), x(k-1)) \quad (1.42)$$

(b) It is obvious.

Exercise 1.7. Commuting functions of a matrix.

Although matrix multiplication does not commute in general

$$AB \neq BA \quad (1.43)$$

multiplication of functions of the same matrix do commute. You may have used the following fact in Exercise 1.5

$$A^{-1} \exp\{At\} = \exp\{At\} A^{-1} \quad (1.44)$$

- (a) Prove that (1.44) is true assuming A has distinct eigenvalues and can therefore be represented as

$$A = Q\Lambda Q^{-1} \quad \Lambda = \begin{bmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_n \end{bmatrix} \quad (1.45)$$

in which Λ is a diagonal matrix containing the eigenvalues of A , and Q is the matrix of eigenvectors such that

$$Aq_i = \lambda_i q_i, \quad i = 1, \dots, n \quad (1.46)$$

in which q_i is the i th column of the matrix Q .

- (b) Prove the more general relationship

$$f(A)g(A) = g(A)f(A) \quad (1.47)$$

in which f and g are any functions definable by Taylor series.

- (c) Prove that (1.47) is true without assuming the eigenvalues are distinct.
Hint: use the Taylor series defining the functions and apply the Cayley-Hamilton theorem [1].

Answer 7. By using Cayley-Hamilton theorem, we know the matrix functions definable by Taylor series can be written in the sum of finite power terms of matrix A [2]. So all the question is answered obviously by matrix commutability of finite power terms.

Exercise 1.8. Finite difference formula and approximating the exponential.

Instead of computing the exact conversion of a continuous time to a discrete time system as in Exercise 1.5, assume instead one simply approximates the time derivative with a first-order finite difference formula

$$\frac{dx}{dt} \approx \frac{x(t_{k-1}) - x(t_k)}{\Delta} \quad (1.48)$$

with step size equal to the sample time, Δ . For this approximation of the continuous time system, compute \tilde{A} and \tilde{B} so that the discrete time system agrees with the approximate continuous time system at the sample times. Comparing these answers to the exact solution, what approximation of $e^{A\Delta}$ results from the finite difference approximation? When is this a good approximation of $e^{A\Delta}$?

Answer 8. In this case, the system function can be written as

$$\frac{x(t_{k+1}) - x(t_k)}{\Delta} = Ax(t_k) + Bu(t_k) \quad (1.49)$$

which could be written as

$$x(t_{k+1}) = (I + \Delta A)x(t_k) + \Delta Bu(t_k) \quad (1.50)$$

So $\tilde{A} = I + \Delta A$, $\tilde{B} = \Delta B$, and $e^{A\Delta} \approx I + \Delta A$. When Δ is small, this is a good approximation.

Exercise 1.9. Mapping eigenvalues of continuous time systems to discrete time systems. Consider the continuous time differential equation and discrete equation

$$\frac{dx}{dt} = Ax \quad (1.51)$$

$$x^+ = \tilde{A}x \quad (1.52)$$

and the transformation

$$\tilde{A} = e^{A\Delta} \quad (1.53)$$

Consider the scalar A case.

- What A represents an integrator in continuous time? What is the corresponding A value for the integrator for discrete time?
- What A give purely oscillatory solutions? What are the corresponding \tilde{A} ?
- For what A is the solution of the ODE stable? Unstable? What are the corresponding \tilde{A} ?
- Sketch and label these A and \tilde{A} regions in two complex-plane diagrams.

Answer 9. (a) In continuous time, $A = 0$, and the corresponding $\tilde{A} = I$.

- A with all eigenvalues on the image axis give purely oscillatory solutions, and the corresponding \tilde{A} has all its eigenvalues on the unit circle.
- A with all eigenvalues on the left-half plane give the solution of the ODE stable (right-half plane give the solution of the ODE unstable), and the corresponding \tilde{A} has its all eigenvalues in the unit circle (out of the unit circle).
- See Fig.1.4. The orange area denote the stable A , and the blue area denote the unstable A .

Exercise 1.10. State space realization

Define a state vector and realize the following models as the state models by hand. One should do a few by hand to understand what the Octave or MATLAB calls are doing. Answer the following questions. What is the connection between the poles of G and the state space description? For what kinds of $G(s)$ does one obtain a nonzero D matrix? What is the order and gain of these systems? Is there a connection between order and the numbers of inputs and outputs?

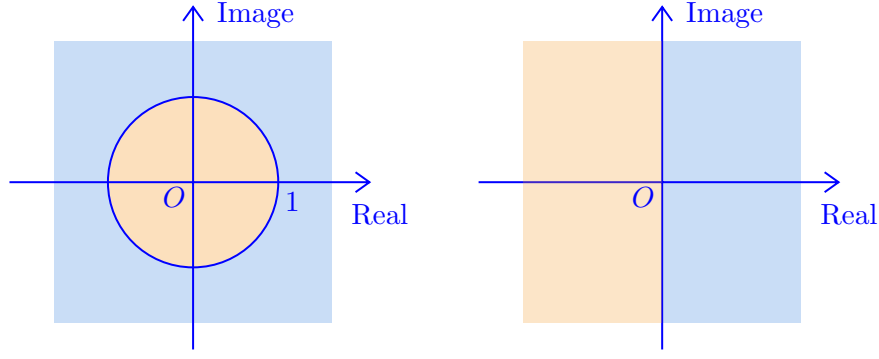


Figure 1.4: Stable and Unstable A and \tilde{A} in two complex plane

(a)

$$G(s) = \frac{1}{2s + 1}$$

(b)

$$G(s) = \frac{1}{(2s + 1)(3s + 1)}$$

(c)

$$G(s) = \frac{2s + 1}{3s + 1}$$

(d)

$$y(k + 1) = y(k) + 2u(k)$$

(e)

$$y(k + 1) = a_1y(k) + a_2y(k - 1) + b_1u(k) + b_2u(k - 1)$$

Answer 10. A state vector is an array of values that combine to represent the state of a system. It is usually composed of the variables describing the system's position and velocity, as well as any other variables that are necessary to describe its behavior.

(a) The state space representation for $G(s) = 1/(2s + 1)$ is:

$$A = -1/2, \quad B = 1, \quad C = 1/2, \quad D = 0 \quad (1.54)$$

(b) The state space representation for $G(s) = 1/(2s+1)(3s+1) = 3/(3s+1) - 2/(2s+1)$ is:

$$A = \begin{bmatrix} -1/2 & 0 \\ 0 & -1/3 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad C = \begin{bmatrix} -1 & 1 \end{bmatrix}, \quad D = 0 \quad (1.55)$$

(c) The state space representation for $G(s) = (2s + 1)/(3s + 1) = 2/3 + 1/3(3s + 1)$ is:

$$A = -1/3, \quad B = 1, \quad C = 1/9, \quad D = 2/3 \quad (1.56)$$

(d) The state space representation for $y(k + 1) = y(k) + 2u(k)$ is:

$$A = 1, \quad B = 2, \quad C = 1, \quad D = 0 \quad (1.57)$$

(e) The state space representation for $y(k+1) = a_1y(k) + a_2y(k-1) + b_1u(k) + b_2u(k-1)$ is [3]:

$$A = \begin{bmatrix} 0 & 1 \\ a_1 & a_2 \end{bmatrix}, \quad B = \begin{bmatrix} b_1 \\ b_2 - a_1b_1 \end{bmatrix}, \quad C = \begin{bmatrix} 1 & 0 \end{bmatrix}, \quad D = 0 \quad (1.58)$$

The connection between the poles of G and the state space description is that the poles of G are the eigenvalues of the matrix A .

For transfer functions that are not strictly proper, there will be a non-zero matrix D .

The dimension of input and output is not directly related to the dimension of system state. For non minimum implementations, the system dimension can be much larger than the input and output dimensions. It is also possible that the dimensions of the system are small, while the dimensions of the input and output are large.

Exercise 1.11. Minimal realization.

Answer 11. In the previous question, the system state space representation was already a minimal realization.

Exercise 1.12. Partitioned matrix inversion lemma.

Answer 12. These formulas can be directly verified by the definition of matrix inversion without complex transformation.

Exercise 1.13. Perturbation to an asymptotically stable linear system.

Answer 13. We know that $\rho(A) < 1$, let $\epsilon = (1 - \rho(A))/2$, there exists matrix norm $\|\cdot\|$ satisfy $\rho(A) < \|A\| + \epsilon$. This implies there exist $\lambda < 1$ such that $\|A\| < \lambda$. So $\|x^+\| \leq \lambda\|x\| + \|B\|\|u\|$, and from the sum formula of proportional sequence, we could get $\|x^\infty\| \leq C\|u\|$, where C is a real constant. So $x \rightarrow 0$ as $u \rightarrow 0$.

Exercise 1.14. Exponential stability of a perturbed linear system.

Answer 14. Without losing generality, we assume $\|x_{k+1}\| \leq \lambda\|x_k\| + C \exp\{-k\}$, where C is a real constant. Then we can get $\|x_k\| = eC (\lambda^k - e^{-k}) / (e\lambda - 1) + \text{constant} \cdot \lambda^{k-1}$. Obviously, the state decrease exponentially to zero.

Exercise 1.15. Are we going forward or backward today?

Answer 15. (a) For simplicity, we use the same notation as forward DP.

$$\bar{x}^0(w) \tag{1.59}$$

$$\tilde{y}^0(w) = \bar{y}^0(\bar{x}^0(w)) \tag{1.60}$$

$$\tilde{z}^0(w) = \bar{z}^0 \bar{y}^0(\bar{x}^0(w)) \tag{1.61}$$

(b) Forward is more efficient, because it does not need to solve function inversion.

Exercise 1.16. Method of Lagrange multipliers.

Answer 16. (a) We know that

$$\frac{\partial L}{\partial x} = Hx + h - D^\top \lambda = 0 \quad (1.62)$$

$$\frac{\partial L}{\partial \lambda} = d - Dx = 0 \quad (1.63)$$

which can be written as

$$\begin{bmatrix} H & -D^\top \\ -D & 0 \end{bmatrix} \begin{bmatrix} x \\ \lambda \end{bmatrix} = - \begin{bmatrix} h \\ d \end{bmatrix} \quad (1.64)$$

(b) It is obvious that $V = (1/2)x^\top Hx = (1/2)x^\top (D^\top \lambda - h) = (1/2)d^\top \lambda$.

Exercise 1.17. Minimizing a constrained, quadratic function.

Answer 17. From previous exercise, we know that

$$\frac{\partial L}{\partial x} = Hx - A^\top \lambda = 0 \quad (1.65)$$

$$\frac{\partial L}{\partial \lambda} = b - Ax = 0 \quad (1.66)$$

So, $x = H^{-1}A^\top \lambda$. Then, $Ax = AH^{-1}A^\top \lambda = b$, $\lambda = (AH^{-1}A^\top)^{-1}b$, $x = H^{-1}A^\top (AH^{-1}A^\top)^{-1}b$, $V = (1/2)x^\top Hx = (1/2)b^\top (AH^{-1}A^\top)^{-1}b$.

Exercise 1.18. Minimizing a partitioned quadratic function.

Answer 18. The result can be obtained by directly substituting the matrix of this exercise into the previous exercise. The two forms of conversion can be performed through Woodbury matrix identity.

Exercise 1.19. Stabilizability and controllability canonical forms.

Answer 19. (a) The controllability matrix is composed of $A^k B$ in the form of column stacking. Notice that the matrix A is partitioned upper triangular matrix, and B has zero rows. So the controllability matrix has zero rows, and the system is not controllable.

(b) We know that

$$x_1(k) = \sum_{i=0}^{k-1} A_{11}^{k-i-1} B_1 \cdot u(i) + \sum_{i=0}^{k-1} A_{12} x_2(i) + A_{11}^k x_1(0) \quad (1.67)$$

Because (A_{11}, B_1) is controllable, $\sum_{i=0}^{k-1} A_{11}^{k-i-1} B_1 \cdot u(i)$ has a full space image, the x_1 can be controlled from any initial state and target state. And $x_2(k) = A_{22}^k x_2(0)$ is not related to u , so x_2 are uncontrollable modes.

(c) Without losing generality, we assume the system is in the form of controllability canonical form. On the one hand, if a system is not stabilizable, because (A_{11}, B_1) is controllable, we can know that A_{22} is unstable. The unstable eigenvalue of A_{22} implies for some $|\lambda| \geq 1$, $\text{rank}[\lambda I - A, B] < n$.

On the other hand, if for some $|\lambda| \geq 1$, $\text{rank}[\lambda I - A, B] < n$, because (A_{11}, B_1) is controllable, all the eigenvalues support $\text{rank}[\lambda I - A_{11}, B_1] = \dim(x_1)$. So $\text{rank}[\lambda I - A_{22}, 0] < \dim(x_2)$, the unstable eigenvalue comes from A_{22} , and the system is not stabilizing.

Exercise 1.20. Regulator stability, stabilizable systems, and semidefinite state penalty.

Answer 20. (a) Without losing generality, we assume the system is in the form of controllability canonical form. Because (A, B) is stabilizable, there exists a sequence of $\dim(x_1)$ inputs that transfers the x_1 to zero. When $k > \dim(x_1)$, we repeat the strategy of u . Notice that the system trajectories start at every $x_1 = 0$ have the scaling of $\|x_2\|$. Because $\|x_2\|$ decrease exponentially, so the objective function $V(x, u) = \sum_{k=0}^{\infty} x_k^T Q x_k + u^T R u$ is finite, which implies the optimization problem is feasible. On the other hand, the solution is unique since $R > 0$ and the objective function is strict convex with u . So the solution of the LQR problem exists and is unique. This implies to that the objective function is non-increasing with time, and we have $x \rightarrow 0$, $u \rightarrow 0$ as $k \rightarrow \infty$. The LQR is stabilizing.

(b) The difference is that $Q > 0$ is replaced by $Q \geq 0$ and (A, Q) detectable. We prove it by contradiction. Assume $x \rightarrow 0$ is false, then there exists ϵ such that for any N , when $k > N$, there exists $\|x_k\| > \epsilon$. We choose N large enough, such that we can ignore sufficient small u and unobservable states. For the observable states x_3 , we expand the system function from time k to time $k + \dim(x_3) = k + m$, then $[y_k, \dots, y_{k+m}] = [Q; QA; \dots; QA^{m-1}]x_k$. Because $[Q; QA; \dots; QA^{m-1}]$ is full rank, so the norm of $[y_k, \dots, y_{k+m}]$ has a positive lower bound. Select N repeatedly, we can always find $\|Qx_k\|$ has a positive lower bound. This contradicts $x_k^T Q x_k \rightarrow 0$. If Q is not positive semidefinite or (A, Q) is not detectable, the state may not convergent to zero.

Exercise 1.21. Time-varying linear quadratic problem.

Answer 21. Following the chapters in the textbook, we get that

$$\begin{aligned} \Pi(k-1) &= Q(k) + A(k)^T \Pi(k) A(k) \\ &\quad - A(k)^T \Pi(k) B(k) (B(k)^T \Pi(k) B(k) + R(k))^{-1} B(k)^T \Pi(k) A(k), \quad k = N, \dots, 1 \\ \Pi(N) &= Q(N), \\ u_k^0(x) &= K(k)x, \quad k = N-1, \dots, 0 \\ K(k) &= - (B(k)^T \Pi(k+1) B(k) + R(k))^{-1} B(k)^T \Pi(k+1) A(k), \quad k = N-1, \dots, 0 \\ V_k^0(x) &= (1/2) x^T \Pi(k) x \end{aligned}$$

Obviously, this problem also can be solved in closed form like the time-invariant case.

Exercise 1.22. Steady-state Riccati equation.

Answer 22. Whether it is stable A , unstable A , singular A , the final result is the same. The results obtained by iteration of Riccati equation and MATLAB function are completely consistent. the eigenvalues value of $A + BK$ is in the unit cycle. And the larger Q is, the closer the eigenvalues value of $A + BK$ is to the origin.

Exercise 1.23. Positive definite Riccati iteration.

Answer 23. by Woodbury matrix identity, we get

$$\begin{aligned}\Pi(k-1) &= Q + A^\top \Pi A - A^\top \Pi B (B^\top \Pi B + R)^{-1} B^\top \Pi A \\ &= Q + A^\top (\Pi - B (B^\top \Pi B + R)^{-1} B^\top \Pi) A \\ &= Q + A^\top (\Pi^{-1} + B R^{-1} B^\top) A > 0\end{aligned}$$

Exercise 1.24. Existence and uniqueness of the solution to constrained least squares.

Answer 24. If this problem has a solution for every b and the solution is unique, we know $\text{rank}(A) = p$. Otherwise, we can choose b to let $Ax = b$ cannot have a solution. And we know

$$\begin{bmatrix} Q \\ A \end{bmatrix} x = \begin{bmatrix} A^\top \lambda \\ b \end{bmatrix} \quad (1.68)$$

If the rank of $[Q; A]$ is less than n , the problem has infinite solutions. On the other hand, the rank condition of A promise that the linear constraints always represent a subspace. And the rank condition of $[Q; A]$ promise there is unique solution.

Exercise 1.25. Rate-of-change penalty.

Answer 25. (a) The Riccati iteration is [4]

$$K(k) = -(R + B^\top \Pi B)^{-1} (B^\top \Pi(t+1)A + M^\top) \quad (1.69)$$

$$\Pi(k-1) = Q + A^\top \Pi A - (A^\top \Pi B + M)(B^\top \Pi B + R)^{-1} (B^\top \Pi A + M^\top) \quad (1.70)$$

(b) The new system is

$$\begin{bmatrix} x(k+1) \\ u(k) \end{bmatrix} = \begin{bmatrix} A & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x(k) \\ u(k-1) \end{bmatrix} + \begin{bmatrix} B \\ I \end{bmatrix} u(k) \quad (1.71)$$

$$\Delta u(k) = u(k) - [0 \ I] \tilde{x}(k) \quad (1.72)$$

The new matrix are

$$x^\top Q x + u^\top R u + \Delta u^\top S \Delta u \quad (1.73)$$

$$= \tilde{x}^\top \begin{bmatrix} Q & 0 \\ 0 & S \end{bmatrix} \tilde{x} + u^\top (R + S) u - 2\tilde{x}^\top \begin{bmatrix} 0 \\ S \end{bmatrix} u \quad (1.74)$$

Exercise 1.26. Existence, uniqueness and stability with the cross term.

Answer 26. The reparameterizing is

$$v = u + R^{-1} M^\top x \quad (1.75)$$

and the new value function is

$$V = (1/2) \sum_{k=0}^{\infty} x(k)^\top (Q - M R^{-1} M^\top) x(k) + v(k)^\top R v(k) \quad (1.76)$$

and the system changes to

$$x(k+1) = (A - B R^{-1} M^\top) x(k) + B v(k) \quad (1.77)$$

the existence, uniqueness and stability conditions are $(A - B R^{-1} M^\top, B)$ stabilizable, $Q - M R^{-1} M^\top \geq 0$, $(A - B R^{-1} M^\top, Q - M R^{-1} M^\top)$ detectable, and $R \geq 0$.

Exercise 1.27. Forecasting and variance increase or decrease.

Answer 27. (a) A counterexample

$$A = \begin{bmatrix} 0 & 1/2 \\ 1/2 & 0 \end{bmatrix}, \quad P = \begin{bmatrix} 1 & 0 \\ 0 & 100 \end{bmatrix} \quad (1.78)$$

(b) A counterexample

$$A = \begin{bmatrix} 0 & 2 \\ 2 & 0 \end{bmatrix}, \quad P = \begin{bmatrix} 1 & 0 \\ 0 & 100 \end{bmatrix} \quad (1.79)$$

(c) A counterexample

$$A = \begin{bmatrix} 0 & 2 \\ 2 & 0 \end{bmatrix}, \quad P = \begin{bmatrix} 1 & 0 \\ 0 & 100 \end{bmatrix} \quad (1.80)$$

Exercise 1.28. Convergence of MHE with zero prior weighting.

Answer 28. In the absence of noise, if the data length is greater than the observability index, the unique solution can be obtained. When the data is insufficient, the solution belong to a linear subspace. When the data reaches the observability index, the initial state can be directly calculated, and the optimization has unique solution.

Exercise 1.29. Symmetry in regulation and estimation.

Answer 29. After tedious calculation, the conclusion shown in the exercise is obvious.

Exercise 1.30. Symmetry in the Riccati iteration.

Answer 30. The conclusion is obvious through variable replacement.

Exercise 1.31. Detectability and observability canonical forms.

Answer 31. This is the dual of answer 19. Only need to transpose the matrix.

Exercise 1.32. Estimator stability and detectable systems.

Answer 32. From the perspective of observability canonical form, the observable state still satisfies Lemma 1.6. And it is obvious that the norm of unobservable state is controlled by the norm of observable, so the estimator is stable. (By the way, it seems that there is no definition of estimator stability.)

Exercise 1.33. Estimator stability and semidefinite state noise penalty.

Answer 33. (a) Because Q is positive semidefinite, $Q = M^\top M$, $M \in \mathbb{R}^{\text{rank}(Q) \times n}$, where M has full rank. So $G = M^\top$, $\tilde{Q} = MM^\top$. Perform a invertible transformation on the system so that only the first $\text{rank}(Q)$ column of B is non-zero, then delete all zero columns of B and delete the corresponding control input. Transform the system back and get $G \in \mathbb{R}^{n \times \text{rank}(Q)}$, $\tilde{Q} = G^\top G \in \mathbb{R}^{\text{rank}(Q) \times \text{rank}(Q)}$.

- (b) The noise is not independent, but linearly correlated.
- (c) This is obvious from the way G is constructed in the first question.
- (d) The difference between this question and Exercise 1.32 is that here (A, G) is stabilizable. For uncontrollable modes, as time increases, it will converge to zero, so the error is bounded. For controllable modes, that does not satisfy the system equation cannot increase the norm of w , so the term w in the objective function is zero only when the system equation is satisfied, so the estimator will not be unstable.
- (e) The estimator will be unstable and the state may diverge.

Exercise 1.34. Calculating mean and variance from data.

Answer 34. (a) We know that

$$\mathbb{E}\hat{x} = \mathbb{E} \frac{1}{N} \sum_k^N x_k = \frac{1}{N} \sum_k^N \mathbb{E}x_k = \mathbb{E}x \quad (1.81)$$

(b) We know that

$$\mathbb{E}\hat{P} = \frac{1}{N} \mathbb{E} \sum_k^N \left(x_k - \frac{1}{N} \sum_k^N x_k \right)^2 \quad (1.82)$$

$$= \frac{1}{N} \mathbb{E} \sum_k^N x_k^2 + \mathbb{E} \sum_k^N \left(\frac{1}{N} \sum_k^N x_k \right)^2 - \frac{2}{N} \mathbb{E} \sum_k^N \left(x_k \frac{1}{N} \sum_k^N x_k \right) \quad (1.83)$$

$$= \frac{N-1}{N} \mathbb{E} \sum_k^N \left(x_k - \frac{1}{N} \sum_k^N x_k \right)^2 \quad (1.84)$$

$$= \frac{N-1}{N} P \quad (1.85)$$

(c) the unbiased estimate is

$$\hat{P} = \frac{N}{N-1} \sum_k^N \left(x_k - \frac{1}{N} \sum_k^N x_k \right)^2 \quad (1.86)$$

(d) $N \geq 100$.

Exercise 1.35. Expected sum of squares.

Answer 35. We focus on one of the quadratic summation terms

$$\mathbb{E}(x_i x_j q_{ij}) = q_{ij} \mathbb{E}(x_i x_j) = q_{ij} (\mathbb{E}(x_i) \mathbb{E}(x_j) + \text{Cov}(x_i, x_j)) = q_{ij} m_i m_j + q_{ij} p_{ij} \quad (1.87)$$

And if we resum, we get what we want to prove.

Exercise 1.36. Normal distribution.

Answer 36. We just have to prove that the standard normal distribution has zero mean and one variance, and the formula we want to prove in this case can be given by a linear transformation of the random variable. The symmetry of the function tells us that the mean is zero, and the variance is calculated by multiplying two identical integrals together and using polar transformations. It can also be calculated using parametric integrals. Obviously, the maximum point is taken at the mean, using what we learned about function monotonicity.

Exercise 1.37. Conditional densities are positive definite.

Answer 37. Since the submatrices of a positive definite matrix are also positive definite, the first two terms hold. By Schur's complement, the last two terms are equivalent to P positive definite.

Exercise 1.38. Expectation and covariance under linear transformations.

Answer 38.

$$\mathbb{E}Cx = C\mathbb{E}x = Cm_x \quad (1.88)$$

$$\text{Cov}(x) = \mathbb{E}Cxx^\top C^\top - (\mathbb{E}Cx)(\mathbb{E}x)^\top = Cm_xm_x^\top C^\top + CP_xC^\top - Cm_xm_x^\top C^\top = CP_xC^\top \quad (1.89)$$

Exercise 1.39. Normal distributions under linear transformations.

Answer 39. See Exercise 1.41.

Exercise 1.40. More on normals and linear transformations.

Answer 40. First we prove that for random vector $X = \mu + BZ \in \mathbb{R}^n$, where B is full rank. And the elements of random vector Z is i.i.d., and is standard normal distribution. let $\Sigma = BB^\top$, then the PDF of X is

$$f(x) = \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)}$$

Because B is full rank, the inverse transformation of $x \rightarrow \mu + Bx$ is $x \rightarrow B^{-1}(x - \mu)$, and the determinant of Jacobian is $|B^{-1}|$. By probability density transformation formula of multivariate random variable

$$g(x) = f\left(\left(B^\top\right)^{-1}(x - \mu)\right) \|B^{-1}\| \quad (1.90)$$

$$= f\left(\left(B^\top\right)^{-1}(x - \mu)\right) |\Sigma|^{-\frac{1}{2}} \quad (1.91)$$

$$= \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)} \quad (1.92)$$

where $f(x) = (2\pi)^{-n/2}e^{-x^\top x}$ is the PDF of Z . And $X \sim N(\mu, BB^\top)$.

If A is a full rank square matrix, $Y = AX + b = ABZ + A\mu + b \sim N(A\mu + b, ABB^\top A^\top)$, its PDF

$$h(x) = \frac{1}{(2\pi)^{n/2} |A\Sigma A^\top|^{1/2}} e^{-\frac{1}{2}(x-A\mu-b)^\top (A\Sigma A^\top)^{-1}(x-A\mu-b)}$$

If A is a row full rank matrix of m rows, Supplement AB with Q orthogonal to the row of AB . And make it a full rank square matrix, Use 0 to supplement lines A and b

$$\overline{AB} = \begin{pmatrix} AB \\ Q \end{pmatrix} \in \mathbb{R}^{n \times n}, \quad \bar{b} = \begin{pmatrix} b \\ 0 \end{pmatrix} \in \mathbb{R}^n, \quad ABQ^\top = 0, \quad \bar{A} = \begin{pmatrix} A \\ Q \end{pmatrix} \in \mathbb{R}^{n \times n}$$

the PDF of $\bar{Y} = \overline{AB}Z + \bar{A}\mu + \bar{b}$ is

$$\bar{h}(x) = \frac{1}{(2\pi)^{n/2} |\overline{AB} (\overline{AB})^\top|^{1/2}} e^{-\frac{1}{2}(x-\bar{A}\mu-\bar{b})^\top (\overline{AB} (\overline{AB})^\top)^{-1}(x-\bar{A}\mu-\bar{b})}$$

So the PDF of Y is

$$\begin{aligned} h(x) &= \int_{x_{m+1}, \dots, x_n} \frac{1}{(2\pi)^{n/2} |\overline{AB} (\overline{AB})^\top|^{1/2}} e^{-\frac{1}{2}(x-\bar{A}\mu-\bar{b})^\top (\overline{AB} (\overline{AB})^\top)^{-1}(x-\bar{A}\mu-\bar{b})} dx \\ &= \frac{1}{(2\pi)^{m/2} |A\Sigma A^\top|^{1/2}} e^{-\frac{1}{2}(x-A\mu-b)^\top (A\Sigma A^\top)^{-1}(x-A\mu-b)} \end{aligned}$$

This is because

$$\int_y \frac{1}{(2\pi)^{(n-m)/2} |QQ^\top|^{1/2}} e^{-\frac{1}{2}y^\top (QQ^\top)^{-1}y} dy = 1, \quad y = [x_{m+1}, \dots, x_n]^\top$$

Further, if the first m rows of A are full rank, and the last k row is 0, then its PDF

$$h(x) = \frac{\delta(x_i = b_i, i = m+1, \dots, n)}{(2\pi)^{m/2} |A\Sigma A^\top|^{1/2}} e^{-\frac{1}{2}(x-A\mu-b)^\top (A\Sigma A^\top)^{-1}(x-A\mu-b)}$$

If the rank of non row full rank matrix A is m , then there is an orthogonal matrix Q so that the first m rows of QA are full rank, and the last k row is 0, that is, the PDF of $Z = QY = QAX + Qb$ is:

$$l(x) = \frac{\delta(x_i = (Qb)_i, i = m+1, \dots, n)}{(2\pi)^{m/2} |QA\Sigma(QA)^\top|^{1/2}} e^{-\frac{1}{2}(x-QA\mu-Qb)^\top (QA\Sigma(QA)^\top)^{-1}(x-QA\mu-Qb)}$$

where $\delta(\cdot)$ is a characteristic function. So the PDF of $Y = Q^{-1}Z$ is

$$\frac{\delta((Qx)_i = (Qb)_i, i = m+1, \dots, n)}{(2\pi)^{m/2} |QA\Sigma(QA)^\top|^{1/2}} e^{-\frac{1}{2}(Qx-QA\mu-Qb)^\top (QA\Sigma(QA)^\top)^{-1}(Qx-QA\mu-Qb)}$$

For a matrix with non full rank, it cannot be directly transformed because it leads to a positive semi definite covariance matrix, which is not allowed.

Exercise 1.41. Signal processing in the good old days – recursive least squares.

Answer 41. The iteration formula is

$$L \leftarrow L + x^\top x, \quad L_0 = P^{-1} + X^\top X \quad (1.93)$$

$$R \leftarrow L + x^\top y, \quad R_0 = X^\top Y + P^{-1}\theta \quad (1.94)$$

$$\hat{\theta} \leftarrow L^{-1}R \quad (1.95)$$

By solving the optimization problem, this iterative method is obviously consistent with the solution with all the measurements.

Exercise 1.42. Least square parameter estimation and Bayesian estimation.

Answer 42. (a) $\hat{\theta} = (X^\top R^{-1}X)^{-1}X^\top R^{-1}(X\theta + e)$, so $\mathbb{E} \hat{\theta} = (X^\top X)^{-1}X^\top(X\theta) = \theta$, and $\text{Cov} \hat{\theta} = (X^\top R^{-1}X)^{-1}X^\top R^{-1}X(X^\top R^{-1}X)^{-1} = \sigma^2(X^\top X)^{-1}$.

(b) $p(\theta|y) = Cp(y|\theta)p(\theta)$, where $p(y|\theta)$ is the PDF of $\mathcal{N}(X\theta, R)$. So this density is also normal too. The mean is $(X^\top X/\sigma^2 + \bar{P}^{-1})^{-1}(X^\top y/\sigma^2 + \bar{P}^{-1}\bar{\theta})$, the variance is $(X^\top X/\sigma^2 + \bar{P}^{-1})^{-1}$.

(c) If $\bar{P}^{-1} \rightarrow 0$, $m \rightarrow (X^\top X)^{-1}X^\top y$, $P \rightarrow \sigma^2(X^\top X)^{-1}$.

(d) see (a). Different from (a), the weight of quadratic form here is R .

(e) see (a).

(f) $p(\theta|y) = Cp(y|\theta)p(\theta)$, where $p(y|\theta)$ is the PDF of $\mathcal{N}(X\theta, R)$. So this density is also normal too. The mean is $(X^\top R^{-1}X + \bar{P}^{-1})^{-1}(X^\top R^{-1}y/\sigma^2 + \bar{P}^{-1}\bar{\theta})$, the variance is $(X^\top R^{-1}X/\sigma^2 + \bar{P}^{-1})^{-1}$.

Exercise 1.43. Least square and minimum variance estimation.

Answer 43. Only need to prove $KRK^\top \geq (X^\top R^{-1}X)^{-1}$, for all $KX = I$. It is obvious.

Exercise 1.44. 123

Bibliography

- [1] James Blake Rawlings, David Q Mayne, and Moritz Diehl. *Model predictive control: theory, computation, and design*, volume 2. Nob Hill Publishing Madison, WI, 2017.
- [2] The cayley-hamilton theorem and the matrix exponential. <http://web.archive.org/web/20041123013742/http://web.mit.edu/2.151/www/Handouts/CayleyHamilton.pdf>. Accessed: 2004-11-23.
- [3] State-space representation. https://en.wikipedia.org/wiki/State-space_representation. Accessed: 2022-12-18.
- [4] Linear quadratic regulator: Discrete-time finite horizon. <https://stanford.edu/class/ee363/lectures/allslides.pdf>. Accessed: 2022-12-18.