

LECTURE NOTES

EE160: Introduction to Control (Fall 2023)

Jiahao Chen^a

^aSIST 1#D206, ShanghaiTech University, China

ARTICLE HISTORY

Compiled October 24, 2024

ABSTRACT

The 14th edition of a textbook is not making it a better material for beginners, so I decided to draft my own lecture notes based off various resources. These lecture notes are freely available and can be downloaded at <https://faculty.sist.shanghaitech.edu.cn/chenjh/courses/>.

Contents

1	Introduction to Control	5
1.1	What is Time?	5
1.2	System and its Block Diagram	5
1.2.1	Three Different Problems that can be Defined by Using a Block Diagram	6
1.3	What is Control?	6
1.3.1	Feedforward Control	7
1.3.2	Proportional Control	7
1.3.3	Model Predictive Control	7
1.3.4	Sliding Mode Control	8
1.3.5	Dynamic Control / Integral Control / Adaptive Control	9
1.3.6	Linear Quadratic Regulator (LQR)	9
1.3.7	Using States as Control Input	9
1.4	Book Recommendations and Other Resources	10
2	Mathematic Model of Linear Systems	11
2.1	Across Variable and Through Variable	11
2.2	Analogue Systems and Analogue Variables [1, Section 2.2]	11
2.2.1	Force-current analogy	11
2.2.2	Force-voltage analogy	12
2.3	System as Excitation and Response	12
2.3.1	Impulse Response	13
2.3.2	System as Operator	14
2.3.3	Signal Analysis	14
2.4	Convolution	15
2.5	Linear System	16

3 Laplace Transform and Transfer Function	16
3.1 Laplace Transform	17
3.2 Transfer Function	18
3.2.1 Relation between Impulse Function and Transfer Function . . .	18
3.2.2 Pole, Zero, and Gain.	18
3.2.3 Strictly Proper	19
3.3 Block Diagram in s -Domain	19
3.4 Signal Flow Diagram and Mason's Signal-Flow Gain Formula* . . .	19
4 Feedback Control System Characteristics	19
4.1 Open Loop and Closed Loop	20
4.2 Practical Implementation of the Inverted Model Controller*	20
4.3 Model of DC Motor	21
4.3.1 Second Order Model of a DC Motor	21
4.3.2 Simplified Model of a DC Motor	21
4.4 Transient Response Comparison between Open and Closed Loop Control	22
4.5 Steady State Error Comparison between Open and Closed Loop Control	23
4.5.1 Steady State Error of a Step Response	23
4.5.2 Steady State Error of a Ramp Response	24
4.6 Foes	24
4.7 Sensitivity Function	24
4.8 Gang of Six	25
4.9 Error Signal Analysis	27
4.10 Disturbance Rejection	27
4.11 Reference Tracking	27
4.12 Noise Attenuation	27
4.13 Sensitivity to Parameter Variation	27
4.14 Q & A	28
5 Feedback Control System Performance	29
5.1 Test signals	29
5.2 Standard Second Order System	29
5.3 Performance Metrics	30
5.4 Dominant Poles	35
5.5 Complex Plane Root Location and the Transient Response	35
5.6 System Types in Terms of Steady State Error	36
5.7 Discussions	36
6 Stability	37
6.1 Stability and Root Locations	37
6.2 Motivation of a Stability Criterion	37
6.3 The Routh-Hurwitz Stability Criterion	38
6.4 Steps to Determine Stability from Polynomial Coefficients	38
6.5 Strength in Determining BIBO Stability	40
6.6 Weakness in Determining Relative Stability	41
6.7 Discussions on Repeated Poles on Imaginary Axis	42
6.7.1 Signal and System are the Same Thing in s -Domain	42
6.7.2 Inverse Laplace Transform	42
6.7.3 Simulation, Modulation, and Envelope	42
6.7.4 Steady State Frequency Response via Zero-Pole Cancellation .	43

7 Steady State Frequency Response	43
7.1 Frequency Response Complex Gain	44
7.2 Transfer Function as Complex Number	45
7.2.1 Trace Example: Simple Real Pole	45
7.2.2 Trace Example: Proportional-Integral (PI) Regulator	46
7.2.3 Trace Example: Proportional Regulator	46
7.3 Bode Plot of Typical Systems	47
7.3.1 Complex Poles/Zeros	47
7.3.2 Real Pole/Zero	49
7.3.3 Pole/Zero at Origin	49
7.3.4 Constant	49
7.3.5 Asymptotic Curves for Sketching Bode Plot By Hand	49
7.4 Frequency Response Measurement	50
7.5 Performance Specifications in Bode Plot	50
7.6 Relative Stability	51
7.6.1 Motivation	51
7.6.2 Phase Margin and Gain Margin	51
7.6.3 Unstable Pole in Open Loop Transfer Function	52
7.6.4 Relative Stability for Closed Loop System	54
7.6.5 RHP Zero in Open Loop Transfer Function	55
7.7 Vector Interpretation of the Frequency Response Complex Gain	55
7.8 Nonminimum Phase System	56
7.9 Discussions	56
7.9.1 Design Example: Tuning of A Nested Loop Control System	56
7.9.2 Frequency Response of The Gang Members	59
7.9.3 Bode Plot of the Closed Loop System	60
8 Stability in Frequency Domain	61
8.1 Motivation and Prerequisite	61
8.2 Nyquist Plot	62
8.2.1 Steps to Plot a Nyquist Plot	64
8.3 Nichols Chart*	66
9 Root Locus Method	67
9.1 Motivation: are these two systems equivalent?	67
9.2 Tuning Arbitrary Parameter Using Root Locus Method	68
9.3 Rules to Plot Root Locus	68
9.4 Time Delay and Padé Approximation	69
9.5 Revisit Nonminimum Phase System	69
9.6 Rhor's Counter Example	70
10 Lead-Lag Compensator	71
10.1 Overview	71
10.2 Phase Lead	72
10.2.1 Lead Compensator Design Using Bode Plot	72
10.2.2 Lead Compensator Design Using Root Locus	72
10.3 Phase Lag	72
10.3.1 Lag Compensator Design Using Bode Plot	72
10.3.2 Lag Compensator Design Using Root Locus	72
10.3.3 Lead-Lag Compensator Design	73

10.4 Two Different Kinds of Steady State Error	73
10.5 Two Different Kinds of Steady State Error	74
10.6 PID Control in a Control Theory Perspective	74
10.6.1 Practical Derivative Control	74
10.6.2 Practical Integral Control	75
10.6.3 Put'em Together	77
11 State Space Model	77
11.1 Conversion from Transfer Function	77
11.1.1 Controllable Canonical Form	77
11.1.2 Diagonal Canonical Form	77
11.1.3 The Duality	77
11.2 Controllability and Observability	77
11.3 Luenberger Observer	77
11.4 Separation Principle	77
11.5 Time Domain Solution of the State Space Model	77
11.6 Reference Input	77
A Which to Use, MATLAB or Python?	78
B Review Math Concepts: Two Kernels	78
B.1 Kernel in Integral Transform	78
B.2 Kernel in Linear Algebra	79
B.2.1 Revisit Linear Map	79
B.2.2 Eigenvalue and Eigenvector	79
B.2.3 Eigenvalue and Eigenfunction	80
B.2.4 Matrix Exponential	81
C Five Ways Solving Ordinary Differential Equations	81
D Zeros and Zero dynamics	81
E Passivity and Stability Margin	81
F Fourier Analysis and Time-Frequency Domain Analysis	82
F.1 Frequency Response and Fourier Analysis	82
F.2 Time Frequency Domain Analysis	82
G Cascaded Loop Control	83
H Pending Proofs	83
I Textbook Errata	83

1. Introduction to Control

In my first lecture on Sept. 26th, 2023, I was suggesting there should be a second episode for the Youtube video entitled “Animation vs. Math” featuring TSC, the sticker man.¹ But this time, we are going to further allow him to bring a new set of tools, including:

$$u, \text{ and } \frac{d}{dt} \quad (1)$$

where u grants TSC the ability to **control**, and $\frac{d}{dt}$ is the magic operator that brings to life some state variable $x \in \mathbb{R}$ describing a **system**, such that x begins to evolve with **time** t .

In this chapter, let's focus on explaining those three elements we have just mentioned: time, system, and control.

1.1. What is Time?

Is time really a thing?

First of all, let's define time. The time is measured in terms of periodic events. The sun rise and sun set, making a day. The SI unit second is defined in terms of the unperturbed transition frequency of the caesium 133 atom.² The positive direction of time elapse is defined in the second law of thermodynamics: “*The entropy of the universe tends to a maximum, or in loose terms, energy spreads out over time.*” The increase of the entropy of an isolated system indicates the direction of time.

But what is entropy?³ From a macroscopic point of view, entropy (denoted by S) changes whenever there is a transfer of heat:

$$\Delta S = \int_{t_0}^{t_1} \frac{-dQ}{T(t)} \quad (2)$$

where $T(t)$ is the temperature when the dissipated (*note the negative sign*) heat dQ is made, and the differential change of heat energy dQ is the work done by the friction force:

$$\begin{aligned} dQ &= \mathbf{F}_{\text{friction}} \cdot d\mathbf{x} \text{ [J]} \\ Q &= \int_{x_0}^{x_1} \mathbf{F}_{\text{friction}} \cdot d\mathbf{x} \text{ [J]} \end{aligned} \quad (3)$$

where $\mathbf{F}_{\text{friction}}$ is the friction force and $\mathbf{x}(t) \in \mathbb{R}^3$ is a trajectory in space. Think, if there is no longer transfer of heat in a universe, does this mean its time stops evolving?

1.2. System and its Block Diagram

“A system is a collection of interconnected parts that form a larger and more complex whole” [2]. It is widely accepted that a diagram of blocks and connecting arrows

¹TSC stands for The Second Coming, the fourth stick figure that was created by Alan Becker.

²<https://en.wikipedia.org/wiki/Second>

³Fun fact: our yearbook is named “ENTROPY”.

is useful for revealing the interconnection of different parts in a complex system. An example block diagram is shown in Fig. 1. We often call the physical system consisting of actuator and process, the **plant**.

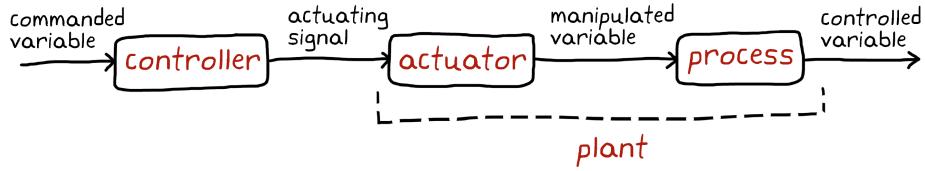


Figure 1. Block diagram showing controller, actuator, and process [2].

Generally speaking, in Fig. 1, a block is often representing an ordinary differential equation (ODE) and an arrow stands for a math operation that is rather bizarre: the convolution. We shall discuss convolution later, and for now you can take it for an integral transform (see Appendix).

1.2.1. Three Different Problems that can be Defined by Using a Block Diagram

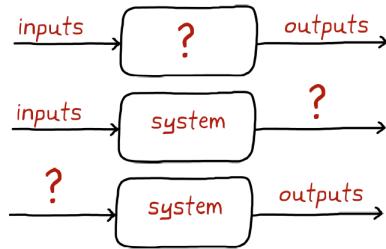


Figure 2. Three different problems arise in control systems [2].

There are three problems that can be studied and they are respectively: system identification problem, simulation problem and control problem [2], depending on which part of the control system is unknown, as shown in Fig. 2.

1.3. What is Control?

The fundamental idea of control is simply we trying to modify the dynamics of any natural process, such that its entropy might decrease or increase at a different pace that is different from natural evolution.

In the class, we have been playing with those math tools in (1) to study how $x(t)$ evolves when its dynamics are one of the followings

$$\frac{d}{dt}x = 1; \quad \frac{d}{dt}x = -1; \quad \frac{d}{dt}x = x; \quad \frac{d}{dt}x = -x; \quad \frac{d}{dt}x = x^2; \quad \frac{d}{dt}x = -x^2 \quad (4)$$

In math, we tend to want to avoid diverging to infinity. The only system among the above that **always** (regardless of initial condition) gives a non-diverging response $x(t)$ is

$$\frac{d}{dt}x = -x \Rightarrow x(t) = x(0)e^{-t}$$

This system evolves with time, and its response $x(t)$ converges towards 0 rather than infinity.

One important goal of control is to make system response to not diverge to infinity. We call a non-diverging system a **stable** system. Stability is a key property of a system and it implies there is an equilibrium point in the system where the dynamics become zero $\frac{d}{dt}x = 0$.

On the contrary, if a system is potential to grow into infinity, the system is **unstable**. For example, $\frac{d}{dt}x = -x^2$ is stable if the initial state satisfies $0 \leq x(0)$, but for other initial state values, the $x(t)$ grows toward infinity as time elapses. We call the range $0 \leq x(0)$ that attracts $x(0)$ towards to a equilibrium point the region of attraction (ROA).

1.3.1. Feedforward Control

To change the dynamics, we need to further append the tool u to any of the former discussed systems, and it yields, e.g.,

$$\frac{d}{dt}x = x^2 + u \quad (5)$$

We want to get rid of the term x^2 by modifying the dynamics of the original system. Assuming x is known, such control goal is simply realized by setting $u = -x^2 + v$, leading to

$$\frac{d}{dt}x = v \quad (6)$$

where term v is yet designed. This means we are essentially treating the term x^2 as a disturbance to the system, and u is able to cancel the effect of such disturbance. Therefore, control is subjective: x^2 is in fact the original dynamics of the system, but it is treated as (internal) disturbance in (5).

1.3.2. Proportional Control

We can further modify the dynamics (6) by designing v to be $v = -x$ to get that nicely behaving system $\frac{d}{dt}x = -x$ again. One might complain that $x(t) = e^{-t}$ is converging too slow. To make the response faster, we can simply let $v = -K_P x$, with $K_P \in \mathbb{R}_+$. This is known as proportional control. Proportional control is the basic form of negative feedback control. **Feedback** refers to the practice to feed the system state (often measured) back to the control input u . **Negative** puts an emphasis that the modified dynamics v must make sure the sign of the exponent of the response $x(t)$ should be negative, e.g., $x(t) = e^{-K_P t}$.

1.3.3. Model Predictive Control

The control input u is applied to an actual system by an actuator. The drawback of the proportional control is that it does not take full ability of the actuator, and the control input is proportional to the system state, requiring the actuator is able to produce an analog signal. In practice, however, the actuator is very likely to operate in an ON-OFF manner. In this case, it makes more sense if we figure out how long the actuator should be turned ON.

For example, let's assume a motor's speed is $x(0) = 10$ rad/s at $t = 0$ s. We are asked to make the motor speed to go up to $x(1) = 100$ rad/s at $t = 1$ s. Let's assume when the actuator is ON, the control input is 100 s^{-2} , which cannot be changed. In this case, the just-exact control input is defined by a Heaviside step function as follows

$$v = 100 [\mathbf{1}(t) - \mathbf{1}(t - 0.9)] \quad (7)$$

The solution to this control system is

$$\begin{aligned} \frac{d}{dt}x &= v \\ \Rightarrow x &= 100 \int_0^t [\mathbf{1}(t) - \mathbf{1}(t - 0.9)] dt \end{aligned} \quad (8)$$

which states that $x(t)$ will ramp up between $[0, 0.9]$ s, as shown in Fig. 3b. Applying

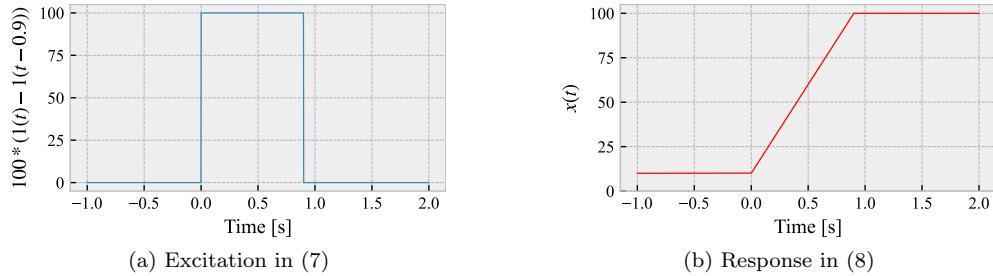


Figure 3. Simple example showing the spirit of model predictive control.

control effort in a period of time just enough to exactly reach the goal describes the spirit of model predictive control.⁴ The exact amount of excitation applied to the system is predicted and calculated based on the model of the system.

1.3.4. Sliding Mode Control

Another idea to apply maximum control effort is the sliding mode control (SMC). In simple terms, the control law is designed to be a signum function as follows

$$v = 100 [\text{s}^{-2}] \operatorname{sgn}(100 [\text{rad/s}] - x)$$

or

$$v = \begin{cases} 100, & \text{if } 100 - x > 0 \\ -100, & \text{others} \end{cases}$$

A similar idea is the bang-bang control or hysteresis control, which sets a region of no control instead of the signum function:

$$v = \begin{cases} 100, & \text{if } 100 - x > 3 \\ -100, & \text{if } 100 - x < -3 \end{cases}$$

⁴Generally, MPC is formulated as optimization problem with inequality constraint.

which gives a non-responsive region having width of 6 rad/s around the goal, 100 rad/s.

1.3.5. Dynamic Control / Integral Control / Adaptive Control

If we view the control input v as the output of some dynamical system: $\frac{dv}{dt} = f(x)$, the resulting control law is called dynamic control.

A well known dynamic control is the integral control, which simply designs v as an integral of the control error:

$$v = K_I \int_0^t (100 - x) dt$$

or in a form of dynamical system:

$$\frac{dv}{dt} = K_I (100 - x)$$

where $K_I \in \mathbb{R}_+$.

Adaptive control also belongs to dynamic control, but with some further considerations for guaranteed stability.

1.3.6. Linear Quadratic Regulator (LQR)

For a simple first-order linear system:

$$\frac{dx}{dt} = ax$$

the LQR law is

$$u = -(a + \sqrt{a^2 + \psi})x$$

with the tuning button $\psi \in \mathbb{R}$. The optimization objective is

$$J_{\text{LQR}} = \int_0^{t_f} [\psi x(t)^2 + u^2] dt$$

where t_f is the final time to end the control. The continuous-time differential Riccati equation will be needed to derive the control law [3, Section 22.4].

1.3.7. Using States as Control Input

For systems of higher order than first-order, using states as control input is a principle that is at the center of control. The idea has been used in (integral) back-stepping control and input-output linearizing control. In addition, the nested loop control or cascaded loop control is designed based on the same principle, but it often assumes the control transients of the inner loop are short enough.

1.4. Book Recommendations and Other Resources

In addition to the books I recommended in the class, please see others' opinions on what control theory is about. Among them, I recommend to watch Brian Douglas's video on Mar. 2nd, 2015 (https://www.youtube.com/watch?v=oBc_BHxw78s): "Why Learn Control Theory". To summarize his video, control theory is widely used in electrical engineering, mechanical engineering, communication engineering, civil engineering, industrial engineering, and aerospace engineering. Control theory is essentially a subject of applied mathematics, it is building models of your systems, it is simulating model to make predictions, it is to understand dynamics and how it interact with environment, it is filtering noises and rejecting disturbances, and it is selecting, building and testing hardware to make sure it has expected performance in an unexpected environment. It is a tool that every engineer should learn to understand his/her system.

2. Mathematic Model of Linear Systems

This chapter, however, does not treat control like what we have done (in time domain) in Chapter 1, as most of those controlled dynamics can only be solved using numerical integrations.

LTI Allowable operations:

Multiply or divide the input by a constant	Integrate or differentiate the input	Add or subtract multiple inputs
$a \cdot x(t)$	$\frac{1}{a} \cdot x(t)$	$\int x(t) dt$ $\frac{d x(t)}{dt}$

$$x_1(t) + x_2(t) \quad x_1(t) - x_2(t)$$

Figure 4. Allowable operations for building a linear time-invariant (LTI) system.

In order to potentially have a closed form solution, we need to study a class of simple systems originated from physics laws. From their governing ODEs, we realize they happen to only consist of the linear operations shown in Fig. 4, hence they are called linear systems. Linear system allows us to apply a series of impulse excitations one by one, and then sum up those impulse responses to produce the output of any arbitrary input, supposing an arbitrary input function can be represented as a series of impulse functions.

2.1. Across Variable and Through Variable

Across variable and through variables are concepts proposed in [1] for helping take abstract of various different physics systems. They are not very useful in this course, and it is sufficient to pay attention to the force-current analogy and force-voltage analogy.

2.2. Analogue Systems and Analogue Variables [1, Section 2.2]

We are going to show that systems that stem from different physics laws end up being very analogue in terms their dynamical equations.

2.2.1. Force-current analogy

The analogy between a damper-spring-mass system and RLC circuit is called force-current analogy.

Kirchhoff's current law states that all currents flowing into a node must be equal to the current flowing out of the node (as a consequence of charge conservation):

$$\frac{v}{R} + C \frac{dv}{dt} + \frac{1}{L} \int_0^t v dt = i(t)$$

where symbols are defined in [1, Fig. 2.3].

From [1, Fig. 2.2]. Newtonian mechanics state that the change of momentum equals

to the sum of forces applied to the particle with mass M :

$$M \frac{d^2}{dt^2}y + b \frac{dy}{dt} + ky = F(t)$$

where M is mass, b is friction/viscous coefficient, k is a spring constant, and F is the force applied to the system. Think which location has been used to define $y = 0$?⁵

2.2.2. Force-voltage analogy

Table 1. Analogy between a particle and a charge.

Particle	Charge
r	Q
$v = \frac{dr}{dt}$	$i = \frac{dQ}{dt}$
m	L
$p = mv$	$\psi = Li$
$F = m \frac{dv}{dt}$	$e = L \frac{di}{dt}$

The analogy between Newton's second law of motion and Faraday's law of induction is called force-voltage analogy, which implies there is an analogy between a particle and a charge, as summarized in Table 1:

- particle's position r and charge (that passes through a cross-sectional area) Q ;
- velocity v and current i ;
- inertial mass m and inductance L ;
- momentum p and flux linkage ψ ;
- force F and voltage e .

Furthermore, the active power in electrical circuit corresponds to the increase in velocity amplitude, and the reactive power corresponds to the change in the direction of the velocity (note velocity should be a vector in space).

2.3. System as Excitation and Response

Alternative to the math equations like ODEs, a system can also be solely defined by its inputs and outputs. This fact actually serves as the foundation of system identification. Input and output are also known as excitation and response.

For the RLC system we just introduced, when the excitation current is described by a Heaviside step function:

$$i(t) = \mathbf{1}(t) = \begin{cases} 1, & t \geq 0 \\ 0, & \text{others} \end{cases} \quad (9)$$

its response is

$$v(t) = K_2 e^{-\alpha_2 t} \cos(\beta_2 t + \phi_2)$$

whose time domain plot is shown in Fig. 5.⁶

⁵It is the equilibrium position of the spring where the force $ky = 0$.

⁶Think why step current excitation produces an impulse response?

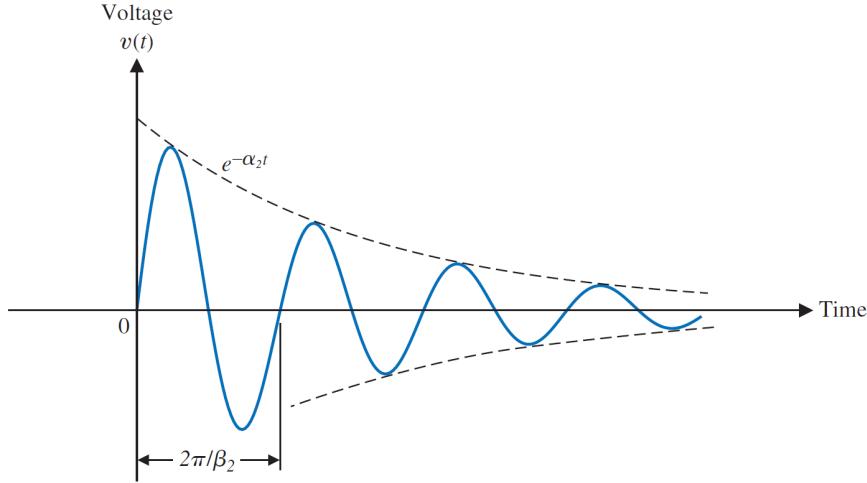


Figure 5. Impulse response of parallel RLC circuit when step current is applied [1, FIGURE 2.3].

For the damper-spring-mass system to have the same response (in waveform), we need to apply an excitation force as follows:

$$F(t) = \frac{d}{dt} \mathbf{1}(t) = \delta(t) \quad (10)$$

which is known as the Dirac delta function. The Dirac delta function is also known as impulse function. It is equivalent to we using a hammer to apply an impact force of 1 N to the system.

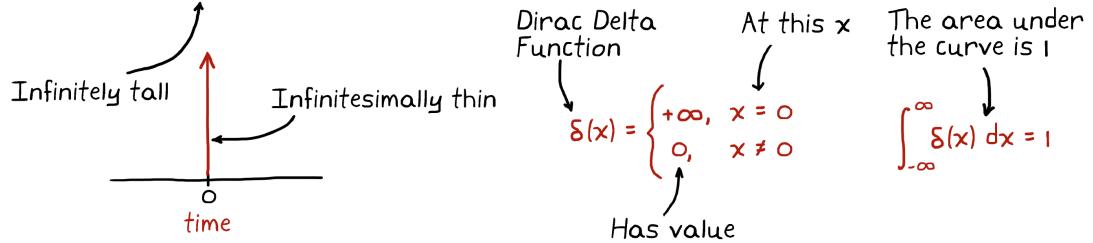


Figure 6. Dirac delta function properties [2].

The Dirac delta function is not a regular function, and some key properties are elaborated in Fig. 6. Recall Heaviside step function is the integral of Dirac delta function, and its derivative of $\mathbf{1}(t)$ at $t = 0$ is infinite. This infinite derivative does not make Heaviside step function to grow into an arbitrary large number, implying such an “impulse” has finite amount of energy.

2.3.1. Impulse Response

Impulse response is the response of an ODE when an impulse function $\delta(t)$ is applied as input. We will see very soon why impulse response is fully representative of an ODE in time-domain. In other words, the ODE and impulse response are equivalent representations of a linear system.

Note even though the response in Fig. 5 is a result of applying a step current in

(9), we still call it the impulse response of this RLC circuit, because the current is the input to a differential-and-integral equation rather than an ODE.

2.3.2. System as Operator

With the concept of impulse response, we are now ready to view system as an operator $f(\cdot)$. In simple terms, the system transfers the impulse excitation $\delta(t)$ into another signal $x(t) = f(\delta(t))$. In general case, when the input is a signal $u(t)$, the system's output becomes $x(t) = f(u(t))$. We are going to show $f(u(t))$ is a convolution of $u(t)$ and $f(\delta(t))$.

An operator should have no memory, otherwise the output will become dependent on the operator's internal states. For example, the impulse response in Fig. 5 would be different, if the capacitor is already charged to some extent at $t = 0$ s. Therefore, to view system as an operator, we need to assume the initial conditions of the system states to be null: $v(0) = \frac{d}{dt}v(0) = 0$, i.e., capacitor voltage and inductance current should be equal to zero when $t = 0$ s.

2.3.3. Signal Analysis

Any signal can be decomposed into a series of delta functions:

$$u(t) \approx \sum_{k=-\infty}^{+\infty} [u(k\Delta t) \Delta t] \delta(t - k\Delta t) \quad (11)$$

with k an integer, where Δt is the sampling period. Motivation is from Calculus: the integral of a signal is the sum of rectangle area under the curve, so we can view any signal $x(t)$ as a bunch of thin rectangles with different heights $x(k\Delta t)$ but the same width Δt . A visualization of (11) is shown in Fig. 7.

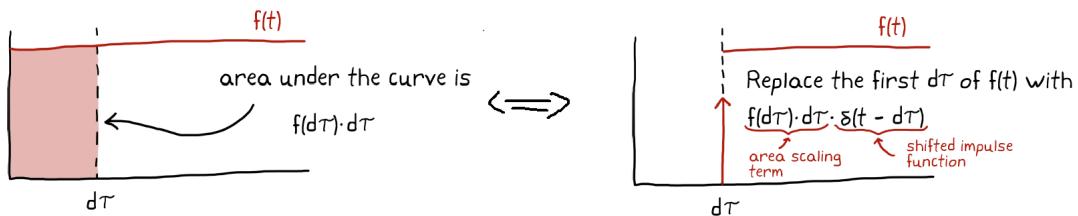


Figure 7. Signal can be equivalent represented using delta function times area under curve.

System with zero initial conditions can be viewed as an operator $f(\cdot)$, so we can apply the system operator to both sides of (11) yields

$$x(t) = f(u(t)) \approx \sum_{k=-\infty}^{+\infty} f(u(k\Delta t)) \delta(t - k\Delta t) \quad (12)$$

which shows that the system response of some arbitrary input function $u(t)$ can be approximately calculated as the sum of a series of impulse responses of delta function excitation of different amplitude of $u(k\Delta t)$.

We need some convenient property for the system operator $f(\cdot)$ so that we can

extract the coefficient outside of the operator $f(\cdot)$ to get following results:

$$x(t) = f(u(t)) \approx \sum_{k=-\infty}^{+\infty} f(u(k\Delta t) \delta(t - k\Delta t)) = \sum_{k=-\infty}^{+\infty} u(k\Delta t) f(\delta(t - k\Delta t))$$

If this holds, we can safely say that a system can be described by its impulse response $f(\delta(t - 0))$. The operator that makes the above result valid is known as the LTI operator, which has been summarized in Fig. 4 before. System's response to any input can be derived as a sum of a series scaled impulse responses, as shown in Fig. 8. This is why impulse response is important because an arbitrary response can be represented as the sum of a series of impulse responses.

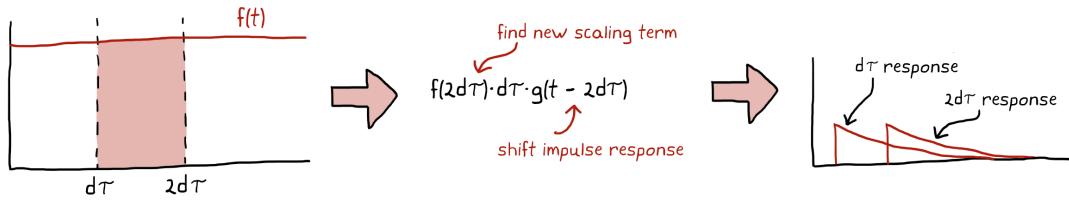


Figure 8. Apply the input slice by slice and sum up the resulting impulse response to get final response, where f is input u in the text, $g(\cdot)$ is impulse response $f(\delta(\cdot))$ in the text, and $d\tau$ is Δt in the text.

2.4. Convolution

Finally, let's formalize the signal decomposition (11) by performing the limit $\Delta t \rightarrow 0$, so sum becomes an integral

$$(11) \xrightarrow{\Delta t \rightarrow 0} u(t) = \int_{-\infty}^{+\infty} u(\tau) \delta(t - \tau) d\tau \triangleq \text{conv}(u(t), \delta(t)) \quad (13)$$

where $k\Delta t$ has been replaced with the variable τ over which the integral (i.e., the sum) is performed, and Δt has been replaced with the differential of time, dt . The integrand in (13) simply means to pick the value of signal $u(t)$ at $t = \tau$.

Now we are ready to define the math operation that an arrow in a block diagram (e.g., Fig. 1) represents. An arrow in a block diagram applies a system's response to another system as input in order to get its response. Therefore, the arrow convolves the previous block's output with the impulse response of the next block.

As another useful property of convolution, a signal $u(t)$ convolves with the delta function $\delta(t - t_1)$ would experience a time shift and becomes $u(t - t_1)$:

$$\text{conv}(u(t), \delta(t - t_1)) = u(t - t_1) \quad (14)$$

In our experiment class, you will find that multiplying two polynomials together can be accomplished by performing the discrete convolution of the polynomial coefficients [2]:

```

1 f = [1 2 3];
2 g = [3 0 1];
3 w = conv(f, g)
4 It prints out: w = 3 6 10 2 3

```

2.5. Linear System

Linear system is a wider concept than linear time-invariant (LTI) system. If a system satisfies homogeneity and superposition, it is a linear system. See Fig. 9 for an intuitive definition.

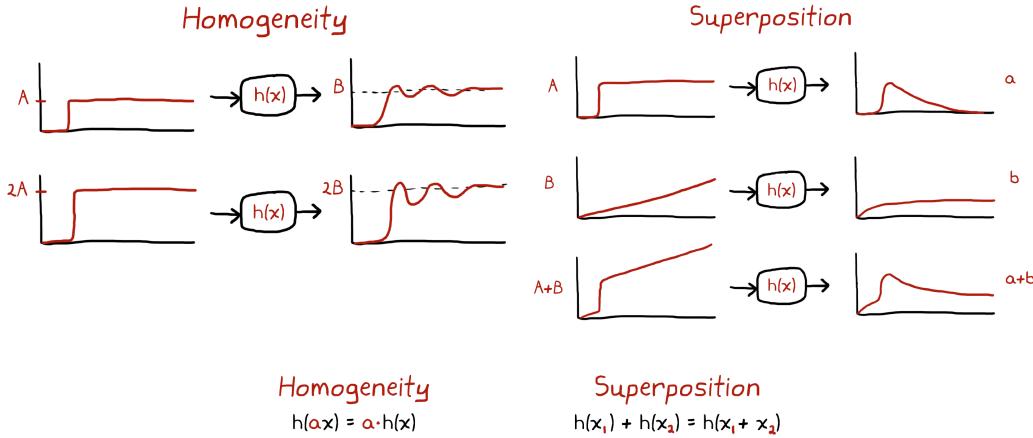


Figure 9. Homogeneity and superposition are two necessary properties of a linear system [2].

If a linear system further sanctifies the property of time-invariance (see Fig. 10), it is then called an LTI system.

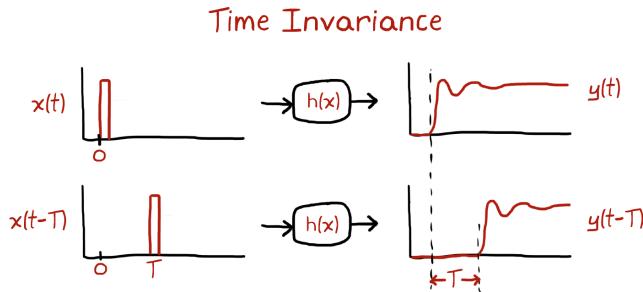


Figure 10. Time invariance states the system should give the same response regardless of the time when excitation is applied [2].

LTI results in sinusoidal fidelity, meaning any sinusoidal signal passing through a system results in a new sinusoidal signal of the same frequency (with a gain in amplitude and a shift in phase). To prove the sinusoidal fidelity, we need to learn to define the frequency response of a transfer function, or in other words, learn solving ODE using Laplace transform, see [1, Chapter 8]

3. Laplace Transform and Transfer Function

Watch Brian Douglas “What are Transfer Functions?”⁷

Those ODEs derived in Section 2.2 are described using LTI operations shown in Fig. 4. We are now ready to have a different representation of those LTI systems,

⁷<https://ww2.mathworks.cn/en/videos/what-are-transfer-functions-1661846920974.html>

which is the transfer function. Transfer function is the output-input ratio after the Laplace transform is applied to an ODE with zero initial conditions. Alternatively, transfer function is defined as the Laplace transform of the impulse response of a system with zero initial conditions.

3.1. Laplace Transform

APPENDIX D Laplace Transform Pairs

$F(s)$	$f(t), t \geq 0$
1. 1	$\delta(t_0)$, unit impulse at $t = t_0$
2. $1/s$	1, unit step
3. $\frac{n!}{s^{n+1}}$	t^n
4. $\frac{1}{(s+a)}$	e^{-at}
5. $\frac{1}{(s+a)^n}$	$\frac{1}{(n-1)!} t^{n-1} e^{-at}$
6. $\frac{a}{s(s+a)}$	$1 - e^{-at}$
7. $\frac{1}{(s+a)(s+b)}$	$\frac{1}{(b-a)} (e^{-at} - e^{-bt})$
8. $\frac{s+\alpha}{(s+a)(s+b)}$	$\frac{1}{(b-a)} [(a-\alpha)e^{-at} - (a-b)e^{-bt}]$
9. $\frac{ab}{s(s+a)(s+b)}$	$1 - \frac{b}{(b-a)} e^{-at} + \frac{a}{(b-a)} e^{-bt}$
10. $\frac{1}{(s+a)(s+b)(s+c)}$	$\frac{e^{-at}}{(b-a)(c-a)} + \frac{e^{-bt}}{(c-a)(a-b)} + \frac{e^{-ct}}{(a-c)(b-c)}$
11. $\frac{s+\alpha}{(s+a)(s+b)(s+c)}$	$\frac{(a-\alpha)e^{-at}}{(b-a)(c-a)} + \frac{(a-b)e^{-bt}}{(c-b)(a-b)} + \frac{(a-c)e^{-ct}}{(a-c)(b-c)}$
12. $\frac{ab(t+\alpha)}{s(s+a)(s+b)}$	$\alpha - \frac{b(\alpha-a)}{(b-a)} e^{-at} + \frac{a(\alpha-b)}{(b-a)} e^{-bt}$
13. $\frac{\omega}{s^2 + \omega^2}$	$\sin \omega t$
14. $\frac{s}{s^2 + \omega^2}$	$\cos \omega t$

Table D.1 continued

D-2 Appendix D Laplace Transform Pairs	
$F(s)$	$f(t), t \geq 0$
15. $\frac{s+\alpha}{s^2 + \omega^2}$	$\frac{\sqrt{\alpha^2 + \omega^2}}{\omega} \sin(\omega t + \phi), \phi = \tan^{-1} \omega/\alpha$
16. $\frac{\omega}{(s+a)^2 + \omega^2}$	$e^{-at} \sin \omega t$
17. $\frac{(s+a)}{(s+a)^2 + \omega^2}$	$e^{-at} \cos \omega t$
18. $\frac{s+\alpha}{(s+a)^2 + \omega^2}$	$\frac{1}{\omega} [(\alpha-a)^2 + \omega^2]^{1/2} e^{-at} \sin(\omega t + \phi),$ $\phi = \tan^{-1} \frac{\omega}{\alpha-a}$
19. $\frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2}$	$\frac{\omega_n}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t} \sin \omega_n \sqrt{1-\zeta^2} t, \zeta < 1$
20. $\frac{1}{s[(s+a)^2 + \omega^2]}$	$\frac{1}{a^2 + \omega^2} + \frac{1}{\omega a^2 + \omega^2} e^{-at} \sin(\omega t - \phi),$ $\phi = \tan^{-1} \frac{\omega}{-a}$
21. $\frac{\omega_n^2}{s[(s^2 + 2\zeta\omega_n s + \omega_n^2)]}$	$1 - \frac{1}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t} \sin(\omega_n \sqrt{1-\zeta^2} t + \phi),$ $\phi = \cos^{-1} \zeta, \zeta < 1$
22. $\frac{(s+\alpha)}{s[(s+a)^2 + \omega^2]}$	$\frac{\alpha}{a^2 + \omega^2} + \frac{1}{\omega} \left[\frac{(\alpha-a)^2 + \omega^2}{a^2 + \omega^2} \right]^{1/2} e^{-at} \sin(\omega t + \phi),$ $\phi = \tan^{-1} \frac{\omega}{\alpha-a} - \tan^{-1} \frac{\omega}{-a}$
23. $\frac{1}{(s+c)[(s+a)^2 + \omega^2]}$	$\frac{e^{-at}}{(c-a)^2 + \omega^2} + \frac{e^{-at} \sin(\omega t + \phi)}{\omega [(c-a)^2 + \omega^2]^{1/2}}, \phi = \tan^{-1} \frac{\omega}{c-a}$

D-1

Figure 11. Screenshot of the Laplace transform pairs from Appendix of [1].

In practice, Laplace transform is as simple as a look-up table, see the screenshot in Fig. 11. The minimum requirement is to remember the Laplace transform of

$$\delta(t), \mathbf{1}(t), t, t^k, e^{-at}, \sin \omega t, \cos \omega t, e^{-at} \sin \omega t, e^{-at} \cos \omega t$$

where k is integer and $a, \omega \in \mathbb{R}$.

Laplace transform can be used to transform an ODE into an algebraic equation with ODE's initial conditions. For example, the damper-spring-mass system is transformed into:

$$\begin{aligned} \mathcal{L} \left[M \frac{d^2}{dt^2} y + b \frac{dy}{dt} + ky \right] &= \mathcal{L}[F(t)] \\ M \left(s^2 Y(s) - sy(0^-) - \frac{dy}{dt}(0^-) \right) + b(sY(s) - y(0^-)) + kY(s) &= F(s) \end{aligned} \quad (15)$$

where $y(0^-)$ and $\frac{dy}{dt}(0^-)$ are called the initial conditions of this second-order ODE. From (15), assuming $y(0^-) = y_0$ and $\frac{dy}{dt}(0^-) = 0$ and solving for $Y(s)$ yield

$$Y(s) = \frac{(Ms+b)y_0 + F(s)}{Ms^2 + bs + k} = \frac{N(s)}{D(s)} \quad (16)$$

which can be transformed into time-domain via inverse Laplace transform. When $F(s) = 0$ and $y_0 \neq 0$, one possible (*depending on values of M, b, k*) impulse response is:

$$y(t) = \mathcal{L}^{-1}[Y(s)] = K_1 e^{-\alpha_1 t} \sin(\beta_1 t + \phi_1)$$

where K_1 , α_1 and β_1 are constants associated with the parameters of the system. When $y_0 = 0$ and $F(t) = \delta(t)$, the solution shares a similar form but has a different initial phase angle than ϕ_1 . Having more than one excitation channels makes the analysis of system response sometimes confusing, and we should stick with one input channel, preferably the input signal $F(t)$.

3.2. Transfer Function

Assuming zero initial conditions, we can derive the ratio between system output and system input for the damper-spring-mass system

$$T(s) = \frac{Y(s)}{F(s)} = \frac{1}{Ms^2 + bs + k} = \frac{N(s)}{D(s)} \quad (17)$$

where $T(s)$ is the transfer function from input $F(s)$ to output $Y(s)$. Neglecting initial conditions, the differential operator s can be defined as follows

$$\begin{aligned} \mathcal{L}\left[\frac{d}{dt}y(t)\right] &= sY(s) - y(0^-) \\ \Rightarrow s &\triangleq \frac{d}{dt}, \text{ if } y(0^-) = 0 \end{aligned}$$

In the sequel, I will always use operator s instead of taking time-derivative $\frac{d}{dt}$. The time-domain and s -domain functions are indicated by its variable, e.g., $\Omega(t)$ and $\Omega(s)$, and sometimes, the lower-case symbol is at the same time replaced with an upper-case symbol, e.g., $y(t)$ and $Y(s)$,

3.2.1. Relation between Impulse Function and Transfer Function

Note the Laplace transform of the impulse function is $1 = \mathcal{L}[\delta(t)]$. Therefore, in s -domain, transfer function is the same as the system's impulse response. In other words, signal and system become the same concept in s -domain.

3.2.2. Pole, Zero, and Gain.

Pole is defined as the s -value that makes a transfer function to become infinity or that makes the denominator polynomial $D(s) = 0$.

Zero is defined as the s -value that makes a transfer function to become zero or that makes the numerator polynomial $N(s) = 0$.

Gain is defined as the transfer function value when $s = 0$ is substituted.

3.2.3. Strictly Proper

Consider a transfer function $T(s) = \frac{\text{Num}(s)}{\text{Den}(s)}$. If the order of the polynomial $\text{Den}(s)$ is equal or higher than that of $\text{Num}(s)$, we say the transfer function is **proper** [4]. In other words, define

$$T(\infty) = \lim_{s \rightarrow \infty} T(s)$$

we have $0 < |T(\infty)| < \infty$ for a proper system, and we have $0 = T(\infty)$ for a strictly proper system.

The strictly proper function can be defined as [4]

$$T_{\text{sp}}(s) = T(s) - T(\infty)$$

3.3. Block Diagram in *s*-Domain

We have mentioned in Section 1.2 that the block in a block diagram is often an ODE, and the arrow in a block diagram is convolution.

The block is a transfer function in *s*-domain, and the arrow between two connected blocks are multiplication in *s*-domain.

A number of blocks in a block diagram can be reduced by applying the [1, Table 2.5]. The fundamental principle, ddd

In my opinion, block diagram has one key advantage over the O.D.E.. In a block diagram, it is not necessary to give a name to all state variables, and it becomes quite easy to pay attention to those state variables that matter.

3.4. Signal Flow Diagram and Mason's Signal-Flow Gain Formula*

Signal flow diagram is only meaningful when the block diagram has too many nodes. In that case, Mason's signal-flow gain formula can be applied to derive the linear dependence between two independent variables in the signal flow graph. While in practical scenarios as far as this course concerns, signal flow graph is the same as block diagram, so it is safe to skip it in this course.

4. Feedback Control System Characteristics

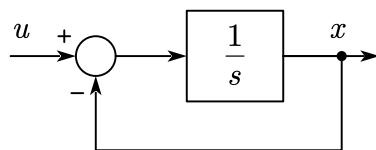


Figure 12. Motivation: the converging response system $sx = -x$ forms a loop.

By making $\frac{1}{s}$ a block in a block diagram, we realize the converging response system $sx = -x$ forms a loop, as shown in Fig. 12. This motivates us that a closed loop might be what we desire for designing a control system that does not have diverging response, which, in most cases, is true.

This section is going to answer why feedback control system is better than a system having no feedback path.

4.1. Open Loop and Closed Loop

Our goal is to make state $x(t)$ follow reference signal $r(t)$.

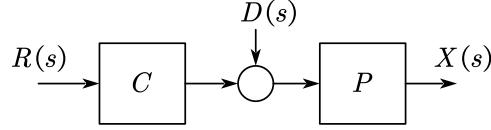


Figure 13. Open loop system.

For an open loop control system like the one in Fig. 13, the full transfer function of the control system is $\frac{X}{R} = CP$. Putting $R(s) = X(s)$ requires $CP = 1$ or $C(s) = P(s)^{-1}$. This kind of controller is known as the inverse system controller, which often is not realizable in practice.

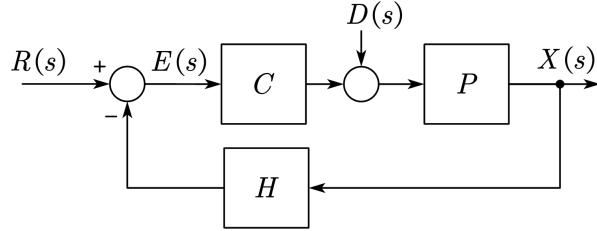


Figure 14. Closed loop system.

As shown in Fig. 14, a closed loop system, on the other hand, gives $\frac{X}{R} = \frac{CP}{1+CPH}$. Note CP is a complex number in nature. As long as $|CP|$ is large enough such that $|CP| \gg 1$, we have $X \approx R$. One realizes that the closed loop control has non-zero error $E(s) \triangleq R(s) - X(s) = \frac{1}{1+CP}R(s)$ in nature, unless $|CP| = \infty$.

4.2. Practical Implementation of the Inverted Model Controller*

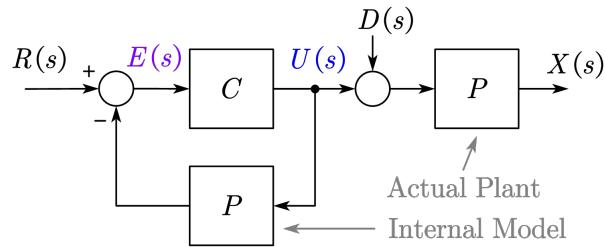


Figure 15. Inverted model control implemented with a feedback loop using an internal model of the plant.

It is interesting to note that a practical way to implement Fig. 13 with $C(s) = P(s)^{-1}$ is to use a feedback loop inside the controller, as shown in Fig. 15 [3]. If we

design $C(s)$ such that the error signal $E(s)$ entering $C(s)$ is close to zero, the inverted model control law is approximately implemented.

Since there is no sensor hardware needed to be applied at the actual plant. The control system is said to be an open loop one. In other words, the controller $C(s)$ has no idea what is going on inside the actual plant. For instance, it is possible the actual plant has been malfunctioning or broken, but there is no way for our controller $C(s)$ to become aware of that fact using an open loop control system. Moreover, there is no way to reject the effect of the disturbance $D(s)$ to the plant, using an open loop control system in Fig. 15.

4.3. Model of DC Motor

Example open loop and closed loop system with a motor can be found in FIGURE 4.12 in [1].

All motor is AC. Even though the voltage applied to the motor can be DC at its terminals, the conductors along the air gap of the motor must carrying an alternating current to maintain a steady torque. In order to provide an alternating current to the conductors, carbon brushes or power electronic devices are necessary to a motor, which realize mechanical and electronic commutation for the current-carrying conductors, respectively.

Assuming perfect conductor commutation, the dc motor consists of a first-order electrical subsystem and a first-order mechanical subsystem. Recall the analogy between the two subsystems is called force-voltage analogy.

4.3.1. Second Order Model of a DC Motor

See FIGURE 4.28 and 4.29 in [1] for a full model of a dc motor when the coil inertia is not neglectable as compared with the rotor inertia. In other words, the disk read head is light in weight. We will address second order plant later.

4.3.2. Simplified Model of a DC Motor

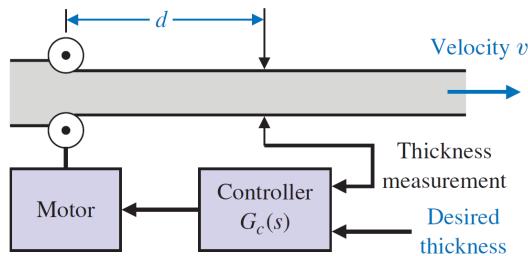


Figure 16. Steel rolling mill control system.

For this chapter, let's consider an simple example of industrial application, and compare the differences when open loop control and closed loop control. The steel rolling mill shown in Fig. 16 [1, FIGURE 4.7] has a heavy rotor, such that the pole of the coil dynamics is far away from the pole of the rotor dynamics, implying that the former can be neglected with limited influence on accuracy. In other words, we can say the mechanical pole dominates the electrical pole. As a consequence, in s -domain,

the dc motor with heavy rotor can be modelled as a first-order transfer function $P(s)$:

$$P(s) = \frac{\Omega(s)}{V(s)} = \frac{g_m}{\tau_m s + 1}$$

where $\Omega(s)$ is the angular speed, $V(s)$ is the voltage applied to the motor armature terminal, g_m is the gain and τ_m is the time constant of the mechanical system.

The block diagram of a simplified DC motor model is shown in [1, FIGURE 4.8], which is a closed loop and the loop reduction $T(s) = X(s)/R(s) = CP/(1 + CPH)$ still works, but it is by definition not a closed loop control system. For a system to be feedback controlled, it has to equip some kind of sensor hardwares.

4.4. Transient Response Comparison between Open and Closed Loop Control

(todo: missing block diagram)

When the expression of a response $x(t)$ has nonzero exponential terms, it is then called a transient response. Closed loop control is able to modify the pole of the closed loop transfer function, so its transient response can be modified to have a larger exponent. To see this, let's compare between the open loop controlled dc motor and the closed loop controlled one.

- The transfer function from reference signal to angular speed using an open loop control is:

$$T(s) = \frac{\Omega(s)}{R(s)} = C(s)P(s) = K_P \frac{g_m}{\tau_m s + 1}$$

which has a real-valued pole $\lambda_1 = -\frac{1}{\tau_m}$.

- The transfer function from reference signal to angular speed using a closed loop control is:

$$T(s) = \frac{\Omega(s)}{R(s)} = \frac{C(s)P(s)}{1 + C(s)P(s)} = \frac{K_P \frac{g_m}{\tau_m s + 1}}{1 + K_P \frac{g_m}{\tau_m s + 1}} = \frac{K_P g_m}{\tau_m s + 1 + K_P g_m}$$

which has a real-valued pole $\lambda_1 = -\frac{1+K_P g_m}{\tau_m}$.

Their time-domain solutions of impulse excitation share the form of

$$\Omega(t) = \mathcal{L}^{-1}\{T(s) \times \mathcal{L}[\delta(t)]\} = K_P g_m e^{\lambda_1 t}$$

which has a larger λ_1 has a faster transient response.

Their s -domain step responses can be derived by substituting $R(s) = 1/s$:

$$\Omega(s) = T(s) \times \mathcal{L}[\mathbf{1}(t)] = \begin{cases} \frac{K_P g_m}{\tau_m s + 1} \frac{1}{s}, & \text{open loop} \\ \frac{K_P g_m}{\tau_m s + 1 + K_P g_m} \frac{1}{s}, & \text{closed loop} \end{cases} = \frac{g_1}{s - \lambda_1} \frac{1}{s}$$

$$g_1 = \frac{K_P g_m}{\tau_m}, \quad \lambda_1 = \begin{cases} \frac{-1}{\tau_m}, & \text{open loop} \\ \frac{-1+K_P g_m}{\tau_m}, & \text{closed loop} \end{cases}$$

For open loop control, we have speed response

$$\begin{aligned}\Omega(s) &= T(s) \times \mathcal{L}[\mathbf{1}(t)] = \frac{g_1}{s - \lambda_1} \frac{1}{s} = \frac{A}{s} - \frac{B}{s - \lambda_1} = \frac{As - A\lambda_1 - sB}{s(s - \lambda_1)} = \frac{g_1}{-\lambda_1} \left(\frac{1}{s} - \frac{1}{s - \lambda_1} \right) \\ &\Rightarrow \begin{cases} A = B \\ -A\lambda_1 = g_1 \end{cases} \Rightarrow A = B = \frac{g_1}{-\lambda_1}\end{aligned}$$

which gives a time-domain response as

$$\Omega(t) = \frac{g_1}{-\lambda_1} \left(1 - e^{\lambda_1 t} \right)$$

See FIGURE 4.13 to have a visualization of the transient response comparison.

4.5. Steady State Error Comparison between Open and Closed Loop Control

With the step response available as $\Omega(t) = \frac{g_1}{-\lambda_1} (1 - e^{\lambda_1 t})$, we can get steady state value by setting $t = \infty$ to get

$$\Omega(\infty) = \frac{g_1}{-\lambda_1} \left(1 - e^{\lambda_1 \infty} \right) = \frac{g_1}{-\lambda_1}$$

which is equivalent to applying final value theorem to the s -domain solution:

$$\Omega(t)|_{t=\infty} = \lim_{s \rightarrow 0} s\Omega(s) = \frac{g_1}{-\lambda_1} s \left(\frac{1}{s} - \frac{1}{s - \lambda_1} \right) = \frac{g_1}{-\lambda_1}$$

4.5.1. Steady State Error of a Step Response

Recall our goal is to make state $x(t) = \Omega(t)$ follow reference signal $r(t)$. It is convenient to evaluate the error signal $e(t) = r(t) - x(t)$ instead. Its s -domain step response is:

$$E(s) = R(s) - \Omega(s) = \frac{1}{s} - \frac{g_1}{-\lambda_1} \left(\frac{1}{s} - \frac{1}{s - \lambda_1} \right) \quad (18)$$

The steady state value of the error signal is

$$e(\infty) = \lim_{s \rightarrow 0} s [R(s) - \Omega(s)] = 1 - \frac{g_1}{-\lambda_1} \quad (19)$$

In order to have zero steady state error, such that $x = \Omega$ coincides with $r(t)$ when t approaches infinity, we need to make $g_1/\lambda_1 = 1$. Think what should the controller $C(s)$ be to make this happen.

For open loop controller, the gain K_P must be tuned to ensure $g_1/\lambda_1 = 1$, assuming the parameters of the system are not time-varying.

For closed loop control, a simple trick to have zero steady state (step) error is to use a infinity loop gain $L(0) = C(0)P(0) = \infty$. This results in an proportional-integral (PI) controller $C(s) = K_P + K_I/s$ with $C(0) = \infty$.

4.5.2. Steady State Error of a Ramp Response

Using the same proportional controller $C(s) = K_P$, the steady state error of a ramp excitation $r(t) = t$ or $R(s) = 1/s^2$ is

$$e(\infty) = \lim_{s \rightarrow 0} sE(s) = s \left[\frac{1}{s^2} - \frac{g_1}{-\lambda_1} \left(\frac{1}{s} - \frac{1}{s - \lambda_1} \right) \right] = \infty \quad (20)$$

which means the proportional control cannot follow a ramping reference signal and its tracking error grows with time.

4.6. Foes

So far, the sole input to our system is the reference signal. Upon analyzing the transient and steady state performance, it seems the best controller design is to use an infinitely large gain in controller gain C . In practice, however, there are at least three input channels to a closed loop control system, at least one of which prevents us from using an infinite gain.

There are undesired phenomena present in a control system, including external disturbance [*measurement noise $n(t)$ and unknown input $d(t)$*] and internal disturbance [*parameter uncertainty ΔP*], leading to degrade in control performance, e.g., causing a remarkable steady state error.

The ultimate goal of the control system design is to keep the reference tracking ability while rejecting all those disturbances to the system. To this end, we need to first introduce the idea of sensitivity function, in order to describe how sensitive to disturbance is our control system.

4.7. Sensitivity Function

The internal disturbance ΔP (which is often a parameter uncertainty) causes a deviation ΔT from T . A metric that evaluates how much perturbation it causes to our system is the sensitivity function, defined by

$$S = \frac{\Delta T(s)/T(s)}{\Delta P(s)/P(s)} \quad (21)$$

where the deviation can be calculated as per definition:

$$\Delta T(s) = \frac{C(P + \Delta P)}{1 + C(P + \Delta P)} - \frac{CP}{1 + CP}$$

In the limit, small incremental changes leads to following definition:

$$S_P^T = \frac{\partial T(s)/T(s)}{\partial P(s)/P(s)} = \frac{\partial \ln T(s)}{\partial \ln P(s)} \quad (22)$$

where the following calculus relation has been substituted:

$$\frac{dx}{x} = d \ln x \Leftrightarrow \int \frac{dx}{x} = \ln x$$

When control system transfer function is $T(s) = \frac{CP}{1+CP}$, the sensitivity function is

$$S_P^T = \frac{1}{1+CP} \quad (23)$$

When control system transfer function is $T(s) = CP$, the sensitivity function is

$$S_P^T = 1 \quad (24)$$

This is the second advantage of using a closed loop control system. The amplitude of the sensitivity function is subject to a factor that is less than 1. Also, it is important to use a **negative** feedback loop, otherwise the denominator in (23) becomes $1 - CP$, making $|S_P^T| > 1$.

In most cases, the transfer function $T(s)$ is a rational fraction:

$$T(s; \alpha) = \frac{N(s; \alpha)}{D(s; \alpha)}$$

where α is a parameter that experiences variation, and N and D are numerator and denominator polynomials in s . As a result, $T(s)$'s sensitivity with respect to parameter α becomes

$$S_\alpha^T = \frac{\partial \ln T}{\partial \ln \alpha} = \left. \frac{\partial \ln N}{\partial \ln \alpha} \right|_{\alpha=\alpha_0} - \left. \frac{\partial \ln D}{\partial \ln \alpha} \right|_{\alpha=\alpha_0} = S_\alpha^N - S_\alpha^D$$

where α_0 is the nominal value of α .

4.8. Gang of Six

Watch video of Douglas “Gang of Six”.⁸

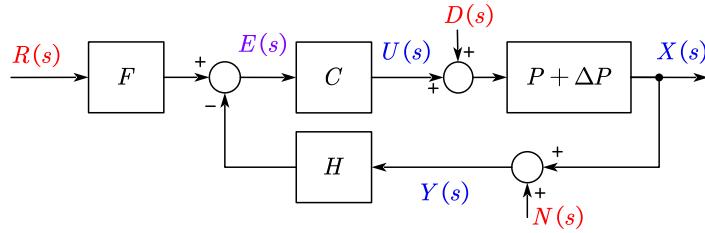


Figure 17. Closed loop system with three input channels.

The closed-loop control system shown in Fig. 17 has considered all three different foes that perturb the control performance. Assuming $H(s) = 1$ in Fig. 17, we can derive

⁸<https://ww2.mathworks.cn/en/videos/control-systems-in-practice-part-8-the-gang-of-six-in-control-theory.html>

the following relationships among the input signals and state/output/input/error:⁹

$$X = \frac{CP}{1+CP}FR + \frac{P}{1+CP}D - \frac{CP}{1+CP}N \quad (25a)$$

$$Y = \frac{CP}{1+CP}FR + \frac{P}{1+CP}D + \frac{1}{1+CP}N \quad (25b)$$

$$U = \frac{C}{1+CP}FR - \frac{CP}{1+CP}D - \frac{C}{1+CP}N \quad (25c)$$

$$E = \frac{1}{1+CP}FR - \frac{1P}{1+CP}D - \frac{1}{1+CP}N \quad (25d)$$

When $F(s) = 1$, the gang of six is reduced as gang of four. We define loop gain as $L \triangleq CP$, and the definitions of the four gang members are now in order:

- Sensitivity function $S = 1/(1 + L)$.
- Complementary sensitivity function is $1 - S$.
- Disturbance sensitivity function is PS .
- Noise sensitivity function is CS .

See also Fig. 18.

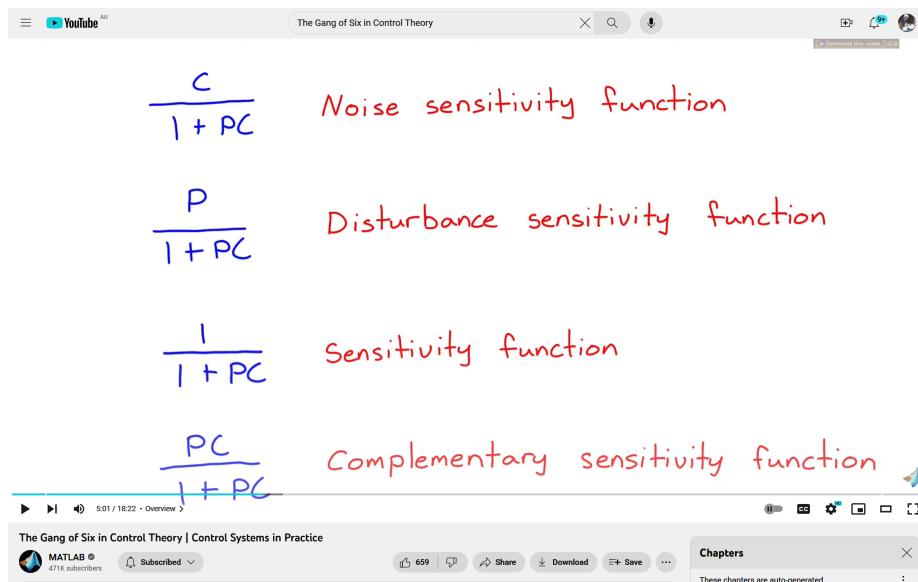


Figure 18. Gang of four by Brian Douglas (Youtube: b8v8scghh8)

⁹In case you wonder how the expression for U is derived (because it appears to be different from the rest):

$$\begin{aligned} U &= CE = C(FR - Y) \\ &= C\left(FR - \frac{CP}{1+CP}FR - \frac{P}{1+CP}D - \frac{1}{1+CP}N\right) \\ &= \left(\frac{C}{1+CP}FR - \frac{CP}{1+CP}D - \frac{C}{1+CP}N\right) \end{aligned}$$

4.9. Error Signal Analysis

Assuming feedforward block $F = 1$, (25d) is rewritten in terms of sensitivity function S as follows:

$$\begin{aligned} E &= \frac{1}{1+L}R - \frac{P}{1+L}D + \frac{L}{1+L}N \\ &= S \times R - PS \times D + (1-S) \times N \end{aligned} \quad (26)$$

4.10. Disturbance Rejection

Using the principle of superposition, let's analyze the effect of external disturbance input D by putting $R = N = 0$:

$$E(s) = -\frac{P}{1+CP}D = -\frac{P}{1+L}D = -PS \times D \quad (27)$$

The disturbance will be rejected if we use a “large” loop gain. Or in rigorous terms, disturbance rejection occurs whenever s is making the gain $|S(s)P(s)|$ small enough.

4.11. Reference Tracking

The error due to change in reference is

$$E(s) = S \times R = \frac{1}{1+L} \times R \quad (28)$$

which suggests that “large” loop gain also minimizes the tracking error.

4.12. Noise Attenuation

The complementary sensitivity function $1 - S$ shows how noise is attenuated in error. Unfortunately, the error excited by noise $N(s)$

$$E(s) = (1 - S) \times N(s) = \frac{L}{1+L} \times N(s) \quad (29)$$

is less attenuated when a “large” loop gain $|L|$ is used. We conclude that there is a compromise between the reference tracking and noise attenuation, because

$$S(s) + (1 - S(s)) \equiv 1 \quad (30)$$

4.13. Sensitivity to Parameter Variation

Uncertainty ΔP affects all three channels of the input. We will take reference input for illustration. Assume $D = N = 0$, and substitute $P + \Delta P$ for P in error analysis

(26) yields

$$\begin{aligned}
 E + \Delta E &= \frac{1}{1 + C(P + \Delta P)} R \\
 \Rightarrow \Delta E &= \left(\frac{1}{1 + C(P + \Delta P)} - \frac{1}{1 + CP} \right) R \\
 &\approx \frac{1}{1 + CP} \frac{\Delta P}{P} R \\
 &= S \frac{\Delta P}{P} R
 \end{aligned} \tag{31}$$

which reveals the reason why S is called as sensitivity function.

4.14. Q & A

Some questions and answers are listed as a brief summary.

Q1: What is a transfer function?

A1: A transfer function is

- a polynomial fraction with a complex variable $s \in \mathbb{C}$;
- an LTI system;
- the system's impulse response in s -domain.

Q2: For a transfer function $T(s) = 1/(s + 2)$, why the time domain response starts from $y(0) = 0.5$ and decays to $y(\infty) = 0$?

5. Feedback Control System Performance

Watch Brian Douglas “The Step Response”¹⁰

This chapter develops performance metric as design requirements for feedback control system design.

5.1. Test signals

Standard test signals include $\delta(t)$, $\mathbf{1}(t)$, t , t^2 , and $\sin \omega t$.

5.2. Standard Second Order System

Closed loop transfer function $T(s)$

$$T(s) = \frac{Y(s)}{R(s)} = \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2} \quad (32)$$

is called the standard second order system, where ω_n is **natural frequency**, and ζ is the **damping ratio**. “Standard” puts an emphasis on the absence of zeros. Can you guess what system’s open loop transfer function is?¹¹

The characteristic equation of transfer function $T(s)$ is its denominator polynomial:

$$s^2 + 2\zeta\omega_n s + \omega_n^2 = 0 \quad (33)$$

whose roots are

$$\begin{aligned} (33) \Rightarrow s &= -\zeta\omega_n \pm \omega_n \sqrt{\zeta^2 - 1} \\ &\triangleq -\tau^{-1} \pm \omega_n \beta \end{aligned} \quad (34)$$

where $\beta \triangleq \sqrt{\zeta^2 - 1}$, and τ is the time constant defined as

$$\tau \triangleq \frac{1}{\zeta\omega_n} \quad (35)$$

which suggests the decaying exponential term in the transient response has an exponent of $-t/\tau$.

The impulse response of a second order system can be obtained using inverse Laplace transform:

$$\mathcal{L}^{-1} \left[\frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2} \right] = \frac{\omega_n}{\sqrt{1 - \zeta^2}} e^{-\zeta\omega_n t} \sin(\omega_n \sqrt{1 - \zeta^2} t), \text{ only when } \zeta < 1 \quad (36)$$

as shown in Fig. 19a. Note the y -axis is re-scaled as y/ω_n and t -axis is re-scaled as $\omega_n t$. Re-scaling the axes will keep the waveform invariant when ω_n changes.

¹⁰<https://ww2.mathworks.cn/en/videos/control-systems-in-practice-part-9-the-step-response-1593067191882.html>

¹¹Hint: divide both numerator and denominator with $s^2 + 2\zeta\omega_n s$.

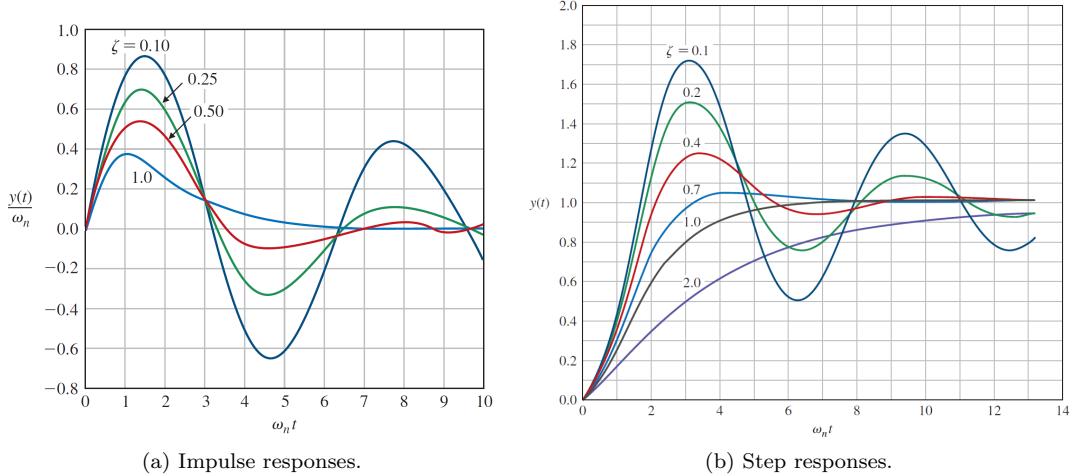


Figure 19. Responses of a standard second order system [1].

The step response of a second order system is

$$\mathcal{L}^{-1} \left[\frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2} \frac{1}{s} \right] = 1 - \frac{1}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t} \sin \left(\omega_n \sqrt{1-\zeta^2} t + \cos^{-1} \zeta \right), \quad \zeta < 1 \quad (37)$$

as shown in Fig. 19b. Note t -axis is re-scaled as ω_{nt} , implying that the waveform shape of the response is not dependent on ω_n . To see this, you can play with the GUI I made with two sliders for adjusting the values of ω_n and ζ , using python packages DearPyGUI and python-control.¹²

From Fig. 19, responses having higher values than final value $y(\infty)$ is said to have an overshoot thus **underdamped**. If the peak of the response is less than the final value $y(\infty)$, then there is no overshoot and the system is said to be **overdamped**. The system is said **critically damped** when the peak value is equal to the final value.

When $\zeta \geq 1$, the second order system is reduced to two first order systems. A general solution can be written as follows: todo

5.3. Performance Metrics

The swiftness of the response is measured by the rise time T_r and the peak time T_p [1]. The 0–100% rise time T_r is mainly used for underdamped system with an overshoot. The 10–90% rise time T_{r1} becomes useful for overdamped system because the time spent in 90–100% would be remarkably long for an overdamped system, which is unfair to be accounted for rising time.

The settling time, T_s , is defined as the time required for the system to settle within a certain percentage δ of the input amplitude (i.e., the command or reference) [1]. In [1], settling time is approximately estimated from

$$e^{-\zeta \omega_n T_s} < 2\% \quad (38)$$

¹²https://github.com/horychen/ee160/blob/master/step_response_visual.py

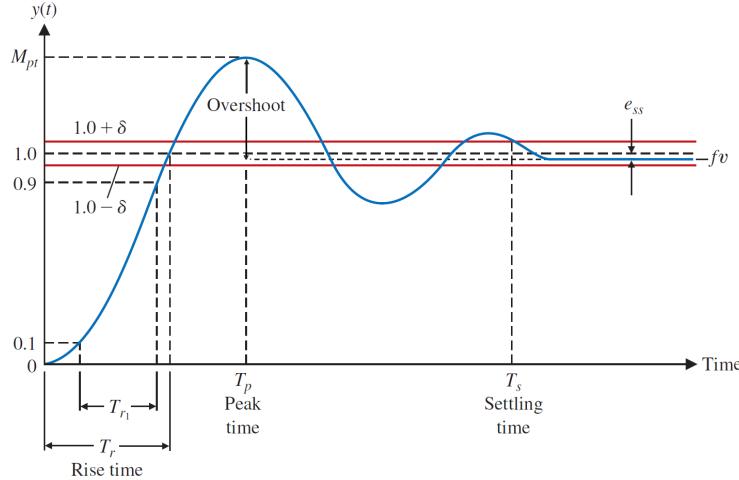


Figure 20. Graphical definitions of the performance metrics in a transient step response.

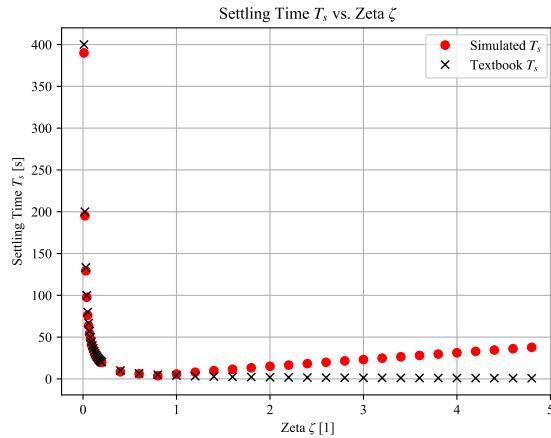


Figure 21. Settling time T_s versus damping ratio ζ .
Premise: step response and standard second order system.

which gives an estimated settling time that is only related to the product of ζ and ω_n :

$$T_s = \frac{4}{\zeta \omega_n} \quad (39)$$

when ζ is not too large. But how large is too large? To have a feel about this, I have used a python snippet to visualize the comparison between the settling time read from graph (results of simulation) and the estimated one $T_s = \frac{4}{\zeta \omega_n}$. The results are shown in Fig. 21. The code snippet is listed here.

```

1 import control
2 from pylab import np, plt, mpl
3 mpl.rc('font', family='Times New Roman', size=10.0)
4 mpl.rcParams['mathtext.fontset'] = 'stix'
5 mpl.rc('legend', fontsize=10)
6 omega_n = 1.0
7 def get_attr(zeta, key='SettlingTime'):
8     T = control.TransferFunction([omega_n ** 2], [1, 2 * zeta *
         omega_n, omega_n ** 2])

```

```

9      return control.step_info(T) ['SettlingTime']
10 zeta_list = np.concatenate(( np.arange(0.0, 0.2, 0.01), np.arange
11     (0.2, 5, 0.2) ))
12 approximated_settling_time_list = [4 / zeta / omega_n for zeta in
13     zeta_list]
14 settling_time_list = [get_attr(zeta, key='SettlingTime') for zeta
15     in zeta_list]
16 plt.plot(zeta_list, settling_time_list, 'o', color='red', label=r'
17     True $T_s$')
18 plt.plot(zeta_list, approximated_settling_time_list, 'x', color='
19     black', label=r'Estimated $T_s$')
20 plt.xlabel(r'Zeta $\zeta$ [1]')
21 plt.ylabel(r'Settling Time $T_s$ [s]')
22 plt.title(r'Settling Time $T_s$ vs. Zeta $\zeta$')
23 plt.grid(); plt.legend(); plt.show()
24 plt.savefig(r'D:\horyc\Desktop\SettlingTimeVsZeta.pdf', dpi=400,
25     bbox_inches='tight', pad_inches=0)

```

From Fig. 21, we learn that the T_s estimate in (39) works quite well when $\zeta < 1$. Therefore, it is safe to say settling time is equal to four times system time constants $4\tau = 4/(\zeta\omega_n)$ [1].

The **steady state error** e_{ss} can be also read on Fig. 20. The response's magnitude at peak time is denoted as M_{pt} . The **percent overshoot** (denoted by P.O.) is defined as

$$\text{P.O.} = \frac{M_{pt} - fv}{fv} \times 100\%$$

where fv is the **final value** of the response. Final value can be calculated as $fv = r(\infty) - e_{ss} = y(\infty)$, where $r(t)$ is reference signal and $e_{ss} = e(\infty)$ is error signal.

For an overdamped system with $\zeta > 1$, theoretically speaking, its peak time can only be read when $t \rightarrow \infty$. When $\zeta < 1$, the peak time can be calculated by differentiating

the response $y(t)$ with respect to time to get¹³

$$\begin{aligned}\frac{d}{dt}y(t) &= \frac{d}{dt} \left[1 - \frac{1}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t} \sin(\omega_n \sqrt{1-\zeta^2} t + \cos^{-1} \zeta) \right] \\ \Rightarrow \frac{d}{dt}y(t) &= \frac{\omega_n}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t} \sin(\omega_n \sqrt{1-\zeta^2} t)\end{aligned}\quad (40)$$

Putting (40) to zero and solving the equation yield an exact expression for peak time

$$\begin{aligned}0 &= \omega_n e^{-\zeta\omega_n t} \frac{\zeta}{\sqrt{1-\zeta^2}} \sin(\omega_n \sqrt{1-\zeta^2} t) \\ \Rightarrow T_p &= \frac{\pi}{\omega_n \sqrt{1-\zeta^2}}\end{aligned}\quad (41)$$

Since this is a theoretical solution, we can in turn use it to validate the accuracy of our numerical simulation. To this end, the previous python script for validating settling time accuracy can be modified to compare between the simulated peak time and the theoretical one in (41). The results are shown in Fig. 22a. From Fig. 22a, the simulated peak time becomes quite off when ζ becomes larger than 0.9. Furthermore, when $\zeta < 0.3$, the zoomed-in plot in Fig. 22b suggests there is remarkable error in simulated peak time—it becomes almost invariant to ζ . This simulation error is probably due to large simulation steps.

¹³The detailed derivation is as follows:

$$\begin{aligned}\frac{d}{dt}y(t) &= \frac{d}{dt} \left[1 - \frac{1}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t} \sin(\omega_n \sqrt{1-\zeta^2} t + \cos^{-1} \zeta) \right] \\ \Rightarrow \frac{d}{dt}y(t) &= - \left(\frac{-\zeta\omega_n}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t} \right) \sin(\omega_n \sqrt{1-\zeta^2} t + \cos^{-1} \zeta) - \frac{1}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t} [\omega_n \sqrt{1-\zeta^2} \cos(\omega_n \sqrt{1-\zeta^2} t + \cos^{-1} \zeta)] \\ \Rightarrow \frac{d}{dt}y(t) &= \frac{\omega_n}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t} \left\{ \zeta \sin(\omega_n \sqrt{1-\zeta^2} t + \cos^{-1} \zeta) - [\sqrt{1-\zeta^2} \cos(\omega_n \sqrt{1-\zeta^2} t + \cos^{-1} \zeta)] \right\} \\ \Rightarrow \frac{\frac{d}{dt}y(t)}{\frac{\omega_n}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t}} &= \zeta \sin(\omega_n \sqrt{1-\zeta^2} t + \cos^{-1} \zeta) - \sqrt{1-\zeta^2} \cos(\omega_n \sqrt{1-\zeta^2} t + \cos^{-1} \zeta) \\ \Rightarrow \frac{\frac{d}{dt}y(t)}{\frac{\omega_n}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t}} &= \sqrt{\zeta^2 + 1 - \zeta^2} \sin \left(\omega_n \sqrt{1-\zeta^2} t + \cos^{-1} \zeta + \arctan \frac{-\sqrt{1-\zeta^2}}{\zeta} \right) \\ \Rightarrow \frac{\frac{d}{dt}y(t)}{\frac{\omega_n}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t}} &= \sin \left(\omega_n \sqrt{1-\zeta^2} t + \cos^{-1} \zeta + \arctan \frac{-\sqrt{1-\zeta^2}}{\zeta} \right) \\ \text{note } \arctan \frac{-\sqrt{1-\zeta^2}}{\zeta} &= -\arccos \frac{\zeta}{1} \\ \Rightarrow \frac{\frac{d}{dt}y(t)}{\frac{\omega_n}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t}} &= \sin \left(\omega_n \sqrt{1-\zeta^2} t + \cos^{-1} \zeta - \arccos \frac{\zeta}{1} \right) \\ \Rightarrow \frac{d}{dt}y(t) &= \frac{\omega_n}{\sqrt{1-\zeta^2}} e^{-\zeta\omega_n t} \sin(\omega_n \sqrt{1-\zeta^2} t)\end{aligned}$$

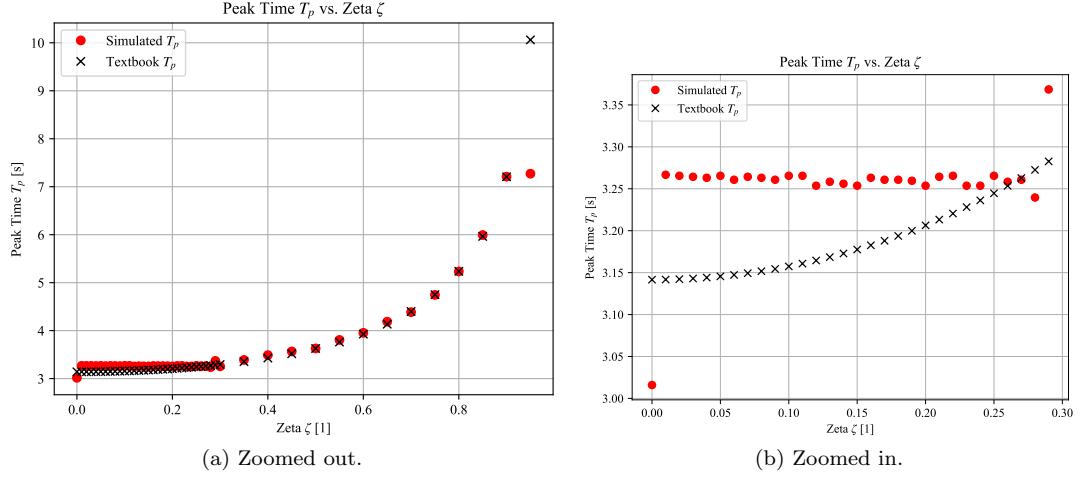


Figure 22. Peak time T_p versus damping ratio ζ . Premise: step response and standard second order system.

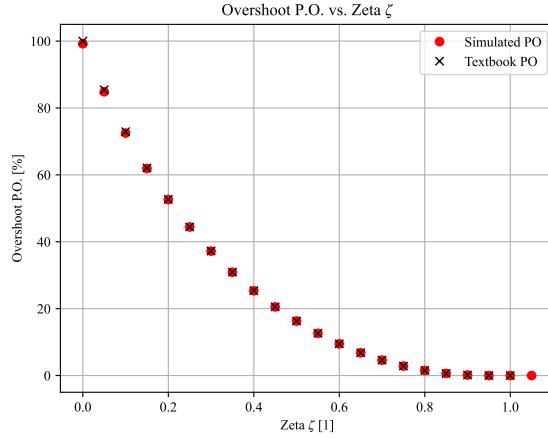


Figure 23. Percentage overshoot P.O. versus damping ratio ζ .
Premise: step response and standard second order system.

The peak response at the estimated peak time (41) is

$$M_{pt} = 1 + e^{-\zeta \frac{\pi}{\sqrt{1-\zeta^2}}}$$

and the resulting estimated percentage overshoot is

$$\text{P.O.} = 100 \times e^{-\zeta \frac{\pi}{\sqrt{1-\zeta^2}}} [\%]$$

which is measured from the command $r(\infty)$ to the peak magnitude M_{pt} , thus there is an minor approximation because the overshoot should be measured from final value $y(\infty)$ to M_{pt} . One can also validate accuracy of P.O.'s expression using the python script, and the results are shown in Fig. 23. Note the overshoot is independent on natural frequency ω_n .

Comparing between Fig. 22a and Fig. 23, one realizes that the textbook peak time T_p and overshoot P.O. are conflicting performance metrics.

Comparing between Fig. 21 and Fig. 23, one finds that settling time T_s and overshoot P.O. both decrease as damping ratio increases. In a physical system, damping ratio

often increases when the value of the dissipator component (e.g., damper and resistor) increases.

As a general design guideline, we need to first pick a ζ value to meet the overshoot requirement. Then the swiftness of the response can be tuned by picking a reasonable ω_n value.

5.4. Dominant Poles

When a system has a large negative real pole, or a pair of conjugate complex poles with large negative real part, it is said those poles far away from the imaginary axis are dominated by other poles that are significantly closer to the imaginary axis. For example, If the far pole $|\lambda_3| \geq 10|\zeta\omega_n|$, then the following approximation is reasonable

$$\frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2} \frac{1}{(s - \lambda_3)} \approx \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2} \quad (42)$$

The roots of its characteristic equation $s^2 + 2\zeta\omega_n s + \omega_n^2 = 0$ are called dominant roots of this third order system.

The concept of dominant poles are valid only when there is no zero near the dominant poles. This can be understood by considering an extreme case in which the zero-pole cancellation occurs. Generally speaking, The poles determine the particular **response modes** (i.e., terms of different exponent) that will be present, and the zeros establish the relative weightings of the individual mode functions. In other words, moving a zero closer to a specific pole will reduce the relative contribution to the output response [1].

5.5. Complex Plane Root Location and the Transient Response

Since the roots of the characteristic equation are complex number, we can mark them in the complex plane (a.k.a. s -plane).

A step response of a general transfer function can be converted into a partial fraction expansion [1, Equation (5.21)], which implies that the transient response consists of terms that have different modes, depending on the location of the characteristic equation roots.

The relation between complex plane root location and the transient response is revealed in [1, FIGURE 5.17]:

- Along the real axis of the s -plane from $-\infty$ to 0, the convergence rate of the response becomes slower and stops to converge when the root is located on the imaginary axis.
- Along the real axis of the s -plane from 0 to ∞ , the diverging speed of the response becomes faster and faster. The response is not bounded by any finite number. Unbounded response is also said to be unstable.
- Along the imaginary axis of the s -plane from $j0$ to $j\infty$, the oscillating frequency of the response becomes higher and higher.

5.6. System Types in Terms of Steady State Error

See Table 5.2 in [1] for the definition of position error constant, velocity error constant and acceleration error constant.

5.7. Discussions

Recall the idea that signal and system are the same thing in s -domain. As a result, the step response of different system does not need to be of similar shape as the step input. On the other hand, a “regular” step response looking response is not necessarily excited by a step input. TODO: add examples.

For higher order systems more than two dominant poles, the performance indices developed in this chapter cannot be directly applied, and the **stability** becomes the highest priority.

6. Stability

A **stable** system is a dynamic system with a bounded response to a bounded input [1]. This input-output property is known as bounded input and bounded output (BIBO) stability [4, Theorem 3.1]. There are other definitions of stability.¹⁴

A stable response is, therefore, a broader concept that includes converging response and bounded response. A typical converging response is $x(t) = e^{-t}$ or $x(t) = 1 - e^{-t}$, and a typical bounded response is $x(t) = 0.707 \sin(t + \pi/4)$. So far, we are only capable of evaluating the stability of a system by looking at its response with respect to a certain input, such as impulse, step or sinusoidal excitation.

Since this chapter, we are going to learn a bunch of tools for analyzing the stability of a control system, including Routh-Hurwitz stability criterion, root locus, Bode plot, Nyquist plot, and Nichols plot. Those tools are especially useful for design of a controller that stabilizes an unstable open loop system. Examples of unstable open loop systems include air-plane, motorcycle, maglev iron ball, inverted pendulum, and bipedal robot.

6.1. Stability and Root Locations

See FIGURE 5.17 for a visualization of the relation between root location and response. There are three stability results depending on the pole locations. [1].

- (1) A system is **stable** if all the poles of the transfer function are in the left hand s -plane (LHP).
- (2) A system is said to have **marginal stability**, if the characteristic equation has **simple roots**¹⁵ on the imaginary axis (e.g., $s = \pm j\omega_u$) with all other roots in the left half-plane. Its steady-state output will be sustained oscillations for a bounded input as long as the input does not contain a frequency that matches the roots on imaginary axis. However, if the input contains a sinusoid (which is bounded) whose frequency ($\omega_u/(2\pi)$) happens to be equal to the magnitude of the imaginary axis roots, then the system response will be unbounded and unstable.
- (3) For an **unstable** system, the characteristic equation has at least one root in the right half of the s -plane or repeated $j\omega$ roots (i.e., non-simple imaginary roots); for this case, the output will become unbounded for any input.

6.2. Motivation of a Stability Criterion

When the characteristic equation, i.e., the denominator of a transfer function, is of low orders, we can resort to the quadratic formula, cubic formula, and quartic formula to get a closed form solution.¹⁶ However, according to Abel–Ruffini theorem¹⁷, there is no closed form solution to polynomial equations of degree five or higher, and an

¹⁴Except BIBO stability, we have, e.g., absolute and relative stability, internal stability, Lyapunov stability, L-2/L-infinity stability, input-to-state stability, asymptotical stability, exponential stability.

¹⁵A simple root is a root with a multiplicity of 1. As a counter-example, we can solve the following large order polynomial equation:

$$(s + 5)^{100} = 0$$

where the root $s = -5$ has a multiplicity of 100, thus not a **simple root**.

¹⁶See, e.g., https://en.wikipedia.org/wiki/Quartic_function#General_formula_for_roots

¹⁷https://en.wikipedia.org/wiki/Abel%E2%80%93Ruffini_theorem

example of non-solvable equation is

$$x^5 - x - 1 = 0 \quad (43)$$

To determine the stability of a system, there is in fact no need to solve for the roots of polynomial, and we only need to determine the signs of the real parts of the roots. For example, for polynomial equation of degree two

$$as^2 + bs + c = 0$$

Vieta's formulas state that the roots λ_1 and λ_2 satisfies

$$\begin{aligned}\lambda_1 + \lambda_2 &= -b/a \\ \lambda_1\lambda_2 &= c/a\end{aligned}$$

As a result, the system must be stable if $-b/a < 0$ and $c/a > 0$, which is a sufficient and necessary condition of stability of this second order polynomial.

For large order polynomial, the Vieta's formulas become only a necessary condition for stability [1, Equation (6.5)] A system is unstable, if the polynomial coefficients do not share the same sign. For example, the following polynomial

$$s^5 + s^4 + s^3 - s^2 + s + 1 = 0$$

is unstable. This motivates that: “sign changing indicates instability”.

6.3. The Routh-Hurwitz Stability Criterion

The generalization of using algebraic combinations of the polynomial coefficients to determining stability, is known as **Routh-Hurwitz stability criterion**.

The best way to learn how to apply this criterion is not to read the determinant based definition in the textbook [1]. Instead, watch the three videos on Routh-Hurwitz stability criterion by Brian Douglas. The screenshot of the first episode is shown in Fig. 24.¹⁸ The table in Fig. 24 is called the Routh array.

The Routh-Hurwitz stability criterion states that

- (1) the number of roots of characteristic equation with positive real parts is equal to the number of changes in sign of the first column of the Routh array;
- (2) and there should be no changes in sign in the first column for a stable system, which is both necessary and sufficient.

6.4. Steps to Determine Stability from Polynomial Coefficients

The full procedure to determine stability from polynomial coefficients are as follows.

S1 Check if all coefficients have the same sign. If not, it is not stable.

S2 Write down Routh array, and if it is a regular case positive and negative numbers in first column, apply the criterion for determining stability.

¹⁸See also <https://www.youtube.com/watch?v=WBCZBOB3LCA>.

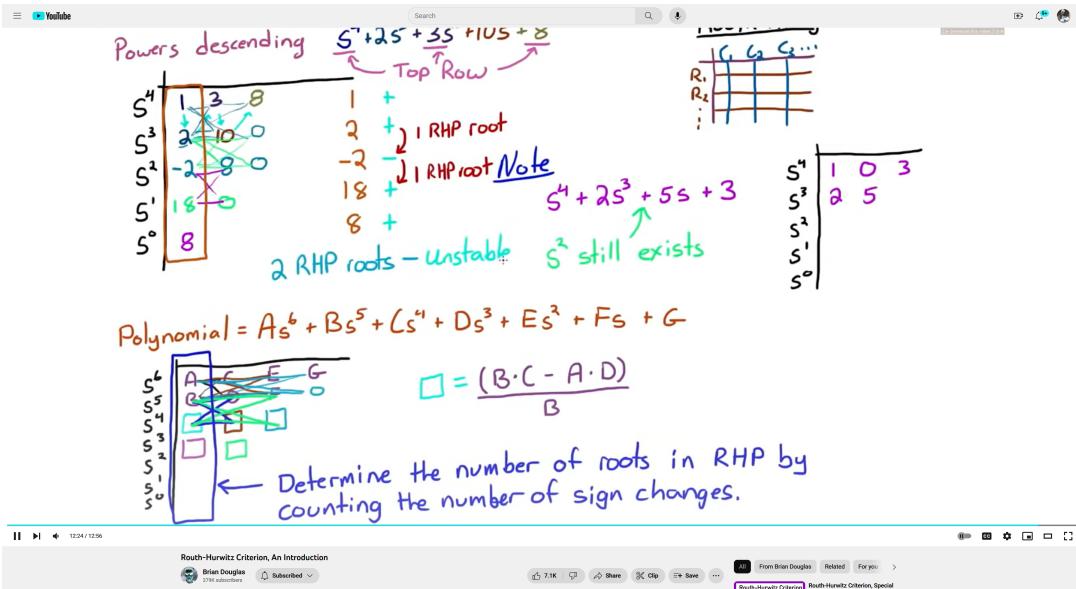


Figure 24. Routh-Hurwitz criterion for regular cases.

S3 If there is a zero in the first column of row s^k of Routh array, we stop. The system is not stable, but we can do more to learn more about the pole locations about the system.

S3.1 Special case one (unstable system): at row s^k , except the first column, there are at least one element being nonzero. We can replace the zeros with a small positive number $\epsilon = 0^+$ at row s^k and continuing to write down Routh array for further information about the number of unstable poles.

S3.2 Special case two (unstable system): at row s^k , including the first column, the rest of elements are all zero. For further information about unstable roots, an auxiliary polynomial is constructed using the row preceding the all-zero row. The order of the auxiliary polynomial is always even and indicates the number of symmetrical root pairs. To be more specific, the auxiliary polynomial contains roots that are symmetrically located about the origin of the s -plane, e.g.,

$$(s + \sigma)(s - \sigma), \\ (s + j\omega)(s - j\omega), \\ \text{or } (s + \sigma + j\omega)(s + \sigma - j\omega)(s - \sigma + j\omega)(s - \sigma - j\omega)$$

The all-zero row is replaced with the coefficients of the time derivative of the auxiliary polynomial. The auxiliary polynomial is a factor of the characteristic polynomial, which can be verified using polynomial division as exemplified in Fig. 25.

A step S2 example is

$$\text{Den}(s) = (s - 1 + j\sqrt{7})(s - 1 - j\sqrt{7})(s + 3) = s^3 + s^2 + 2s + 24$$

$6s^4 + 12s^3 = 0 \Rightarrow s^2 + 2 = 0 = P(s)$

Take derivative of $P(s) = \frac{dP(s)}{ds} = 2s$ 2s replace all zero row with this

One step further...
If you have a $P(s)$
then it is a factor
of the original polynomial.

$P(s) \cdot R(s) = Q(s)$

$Q(s) = (s^2 + 2)(s^3 + 2s^2 + 4s + 6)$

$R(s)$

$\begin{array}{r} s^3 + 2s^2 + 4s + 6 \\ \hline s^2 + 2 \quad | \quad s^5 + 2s^4 + 6s^3 + 10s^2 + 8s + 12 \\ \underline{s^5 + 0s^4 + 2s^3} \\ \hline 2s^4 + 4s^3 + 10s^2 \\ \underline{2s^4 + 0s^3 + 4s^2} \\ \hline 4s^3 + 6s^2 + 8s \\ \underline{4s^3 + 0s^2 + 8s} \\ \hline 6s^2 + 0s + 12 \\ \underline{6s^2 + 0s + 12} \\ \hline 0 \end{array}$

No remainder $\rightarrow 0$

Figure 25. Special case two, auxiliary polynomial and polynomial division.

A step S3 or S3.1 example is

$$\text{Den}(s) = s^5 + 2s^4 + 2s^3 + 4s^2 + 11s + 10$$

A step S3 or S3.2 example is (when $K = 8$)

$$\text{Den}(s) = s^3 + 2s^2 + 4s + K$$

The resulting auxiliary polynomial is $A(s) = 2s^2 + K$.

Another step S3 or S3.2 example is

$$\text{Den}(s) = (s+1)(s+j)(s-j)(s+j)(s-j) = s^5 + s^4 + 2s^3 + 2s^2 + s + 1$$

in which special case two occurs twice. Writing down Routh array will end up with two auxiliary polynomials that are both factors of $\text{Den}(s)$:

- one is $A_1(s) = (s+j)(s-j)(s+j)(s-j)$ at row s^4 of Routh array,
- and the other is $A_2(s) = (s+j)(s-j)$ at row s^2 of Routh array.

which if there is suggests that a valid Routh array

Can you think of a counter-example to reject “the order of the auxiliary polynomial is always even and indicates the number of symmetrical root pairs”? There is none, unless we put the s^0 term to zero, e.g.,

$$\text{Den}(s) = s(s+1)(s+j)(s-j)(s+j)(s-j) = s^6 + s^5 + 2s^4 + 2s^3 + s^2 + s + 0$$

6.5. Strength in Determining BIBO Stability

The advantages using a Routh–Hurwitz criterion is when the transfer function has a polynomial in its denominator instead of pole explicit form. For example, consider the following open loop system:

$$CP(s) = \frac{K}{s^4 + 10s^3 + 35s^2 + 50s + 24} \quad (44)$$

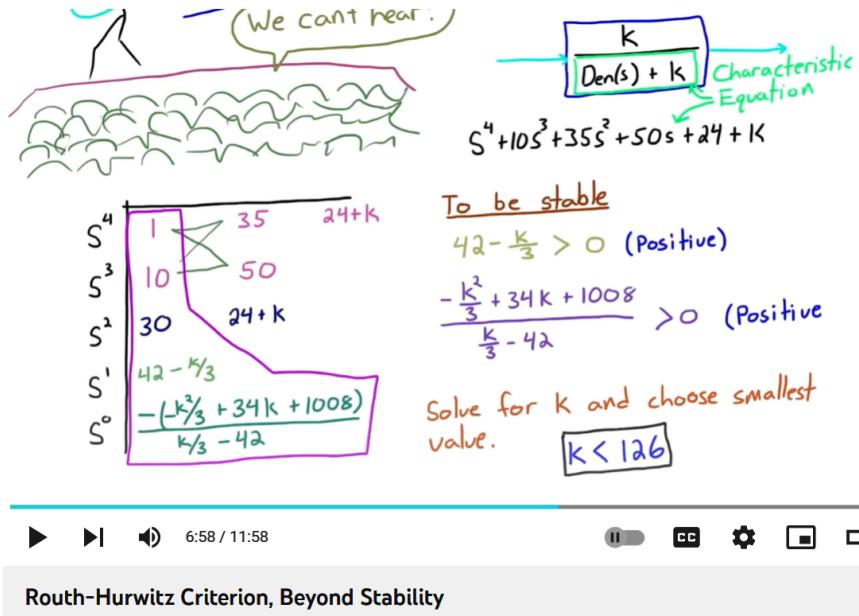


Figure 26. Routh-Hurwitz criterion for controller design of a transfer function having implicit poles.

RH criterion can be used for controller design to ensure the BIBO stability, also known as absolute stability (in contrast to relative stability, soon to be introduced).

In fact, even when the open loop transfer function is provided in its pole explicit form, the Routh array of the characteristic equation of the closed loop system is completely different. Consider the welding control example from textbook, where we can easily find the range of value K such that the following characteristic equation has stable roots:

$$\text{Den}(s) = 1 + C(s)P(s) = 1 + \frac{K(s+a)}{s+1} \frac{1}{s(s+2)(s+3)} \quad (45)$$

where a is a constant.

6.6. Weakness in Determining Relative Stability

In order to have an estimation of how far our roots are away from imaginary axis, we need to substitute, e.g., $s' = s - \sigma$ into the polynomial with $\sigma \in \mathbb{R}$. If the resulting Routh array in terms of s' is stable, we then say that the roots are **relatively stable** with a margin of at least σ away from the imaginary axis. A characteristic root that is farther away from the imaginary axis, corresponds to an exponential term that decays faster, and is said to be **more stable**.

The issue of using RH criterion for determining relative stability, is that we need to determine a maximal value of σ in an iterative fashion, by guessing the value of σ that makes the Routh array in terms of s' marginally stable.

6.7. Discussions on Repeated Poles on Imaginary Axis

Consider the following unstable system

$$Y(s)/U(s) = \frac{1}{(s^2 + 16)^2} \quad (46)$$

with repeated pole $4j$ on the imaginary axis. Since $4j$ is not a simple root, the system is unstable. We will now discuss from 4 different aspects as follows.

6.7.1. Signal and System are the Same Thing in s -Domain

Since signal and system are the same in s -domain, we can switch our perspective a bit. Recall the gain for frequency response to $\sin(4t)$ is infinity for the following marginally stable system with simple roots on imaginary axis:

$$\frac{1}{((4j)^2 + 16)} = \frac{1}{0} = \infty$$

So its frequency response to a sinusoidal input $\frac{1}{((4j)^2 + 16)}$ will be unbounded.

6.7.2. Inverse Laplace Transform

Consider the following inverse Laplace transform: (todo: need to double check)

$$\mathcal{L}^{-1} \left[\frac{4s}{(s^2 + 16)^2} \right] = t \sin(4t)$$

and we conclude the system's impulse response is unbounded, as $t \rightarrow \infty$.

6.7.3. Simulation, Modulation, and Envelope

The following MATLAB snippet numerically simulates the impulse, step and frequency responses of (46).

```

1 %%%
2 close all; cla; clc
3 s = zpk(0, [], 1)
4 % s = zpk('s')
5
6 P = 10*16^2 / ((s+10)*(s^2+16)^2)
7 % P = 10*16^2 / ((s+10)*(s^2+16))
8 subplot(311); step(P)
9 subplot(312); impulse(P)
10
11 t = linspace(0, 10000, 10000);
12 u = sin(3*t);
13 y = lsim(P, u, t);
14 subplot(313); plot(t,y)

```

The frequency response to $\sin(3t)$ is a modulation of two sinusoids of the natural frequency 4 rad/s and forcing frequency 3 rad/s, such that the envelop of the response oscillates at a frequency of 1 rad/s. When the exciting angular speed approaches 4 rad/s, the period of this envelop becomes infinite, hence the response is unbounded.

6.7.4. Steady State Frequency Response via Zero-Pole Cancellation

The gain for frequency response to $\sin(4t)$ is

$$\frac{1}{((4j)^2 + 16)^2} = \frac{1}{0} = \infty$$

The gain for frequency response to $\sin(t)$ is

$$\frac{1}{((1j)^2 + 16)^2} = \frac{1}{15^2}$$

which is finite. In order to produce a steady state frequency response for this system (46), we need to carefully set the initial states so we can avoid any transient response. In other words, the transient response of (46) is always unstable. To understand how it is possible to have a bounded response to (46), let's convert it back to its ODE form

$$(46) \Rightarrow (s^2 + 16)^2 y = u \Rightarrow s^4 y + 32y^2 + 16^2 y = u$$

where s has been used as the differential operator d/dt . Recall that when we apply a Laplace transform to above ODE, we have:

$$\left(s^4 Y - \frac{d^3 y}{dt^3}(0^-) - s \frac{d^2 y}{dt^2}(0^-) - s^2 \frac{dy}{dt}(0^-) - s^3 y(0^-) \right) + 32 \left(s^2 Y - s \frac{dy}{dt}(0^-) - s^2 y(0^-) \right) + 16^2 Y = U$$

With proper initial conditions and no input ($U(s) = 0$), it is possible to introduce zeros to the expressions of $Y(s)$ such that proper pole zero cancellation occurs in the fraction, thus a bounded response becomes possible.

All these efforts of Laplace transform considering initial conditions can be evaluated by setting $s = 0 + j\omega$ in frequency response analysis. We shall address the complex gain of the sinusoidal response in next chapter.

7. Steady State Frequency Response

In the discussions at the end of last chapter, the frequency response has been used to find the exact exciting frequency that leads to an unbounded output of a marginally stable system. However, bear in mind that the steady state gain of the frequency response does not always reveal the instability of the system, e.g., the system with repeated poles on $j\omega$ -axis, which is a con of the frequency response analysis. In fact, if there is no transient at all, it is even possible to make unstable system $1/(s^2 + 16)^2$ to have a bounded response (with no input but nonzero initial states). On the other hand, if you try to draw a response due to a sinusoidal input with zero initial states, a phase shift between the input and the response is built upon the transients—the response is no perfect sinusoid when $t \in [0, \infty]$ (given the fact that both the input and the response start from zero). A simple example is the integrator $1/s$, for which when the input is a sine, the output experiences transients and converges to a cosine; but when the input is a cosine, the output is a perfect sinusoid with no transients.

7.1. Frequency Response Complex Gain

Previously, we have evaluated the steady state response of step, ramp, and parabolic input, in terms of its steady state error. Let's further extend the concept of steady state evaluation of the response to arbitrary sinusoidal inputs $R(s) = \frac{\omega}{s^2 + \omega^2}$.

Frequency response is the system's steady state response to sinusoidal inputs, in which the transients are not important thus shall be neglected. To this end, replacing $s = \sigma + j\omega$ with $j\omega$ in $T(s)R(s)$ provides steady state frequency response. Alternatively, we can prove above " $\sigma + j\omega \rightarrow j\omega$ " trick by considering the following example. Let the excitation be $r(t) = \sin \omega t$ or $R(s) = \frac{\omega}{s^2 + \omega^2}$, a second order system's response is

$$Y(s) = T(s)R(s) = \frac{N(s)}{(s - \lambda_1)(s - \lambda_2)} \frac{\omega}{s^2 + \omega^2} = \frac{g_1}{s - \lambda_1} + \frac{g_2}{s - \lambda_2} + \frac{\alpha s + \beta \omega}{s^2 + \omega^2} \quad (47)$$

where α and β are coefficients to be determined from the partial fraction expansion and are related to coefficients of $N(s)$ and poles λ_1, λ_2 . The time-domain response is

$$\begin{aligned} y(t) &= g_1 e^{\lambda_1 t} + g_2 e^{\lambda_2 t} + \mathcal{L}^{-1} \left[\frac{\alpha s + \beta \omega}{s^2 + \omega^2} \right] \\ &= g_1 e^{\lambda_1 t} + g_2 e^{\lambda_2 t} + \alpha \cos \omega t + \beta \sin \omega t \end{aligned}$$

In the limit $t \rightarrow \infty$, the first two exponential terms vanish, and we have¹⁹

$$\begin{aligned} y(t)|_{t \rightarrow \infty} &= \mathcal{L}^{-1} \left[\frac{\alpha s + \beta \omega}{s^2 + \omega^2} \right] = \alpha \cos \omega t + \beta \sin \omega t \\ &= \text{sgn}(\alpha) \sqrt{\alpha^2 + \beta^2} \cos \left[\omega t + \arctan \frac{-\beta}{\alpha} \right] \\ &= |T(j\omega)| \sin [\omega t + \angle T(j\omega)] \\ &\Leftarrow \begin{cases} |T(j\omega)| = \text{sgn}(\alpha) \sqrt{\alpha^2 + \beta^2} \\ \angle T(j\omega) = \arctan \frac{-\beta}{\alpha} \end{cases} \end{aligned}$$

where $\alpha \neq 0$ and if the amplitude is negative (when $\alpha < 0$), it should be accounted for phase shift of 180° . It is, therefore, suggested to use arctan2 instead of \arctan .

Final value theorem cannot be used to attain $y(t)$ as $t \rightarrow \infty$, because final value theorem can only be applied to a response $Y(s)$ when $Y(j\omega)$ exists with $j\omega \neq 0$. A non-rigorous proof to show that the last row of equation holds is as follows.

$$\begin{aligned} Y(s) &= T(s)R(s) = T(s) \frac{\omega}{s^2 + \omega^2} = 0 + 0 + \frac{\alpha s + \beta \omega}{s^2 + \omega^2} \\ \Rightarrow T(s) &= \frac{\frac{\alpha s + \beta \omega}{s^2 + \omega^2}}{\frac{\omega}{s^2 + \omega^2}} = \frac{\alpha s + \beta \omega}{\omega} \\ \Rightarrow T(j\omega) &= \frac{\alpha j\omega + \beta \omega}{\omega} = \frac{\alpha j + \beta}{1} \\ \Rightarrow |T(j\omega)|^2 &= \alpha^2 + \beta^2 \end{aligned}$$

¹⁹Trigonometry identity used here can be found at https://en.wikipedia.org/wiki/List_of_trigonometric_identities#Sine_and_cosine

In the sequel, whenever we substitute $s = j\omega$, it assumes the above process has been performed, including applying a sinusoidal excitation of the angular speed ω , checking whether or not the system is stable, neglecting exponential decaying terms, and focusing on the amplitude gain $|T(j\omega)|$ and phase shift $\angle T(j\omega)$. In other words, the transfer function $T(s)$ becomes a **steady state complex gain** $T(j\omega)$ that is function of the frequency of a sinusoidal input.

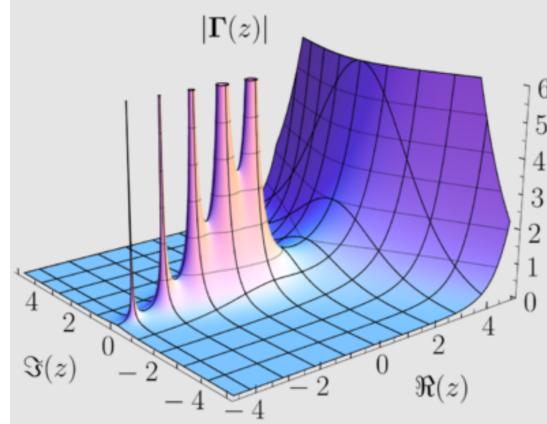


Figure 27. Example of complex-valued function.

As a side note here, the transfer function is only engineering valid and is not theoretically sound, and the step assuming $s = j\omega$ makes it not reliable for determining the stability of the system. However, $s = j\omega$ does allow us to simplify the analysis of a complex-valued function, see e.g., Fig. ??.²⁰

7.2. Transfer Function as Complex Number

The transfer function, the function of complex number $s = 0+j\omega$, is a complex number. Therefore, we can draw a trace of the transfer function in the complex s -plane as ω varies.

7.2.1. Trace Example: Simple Real Pole

See FIGURE 8.3 for the transfer function trace of a RC circuit having time constant of $1/\omega_1$. We can prove that the trace is a segment of a circle by checking the vectors pointing at a point (x, y) from $(0, 0)$ (when $\omega = \infty$) and $(1, 0)$ (when $\omega = 0$) on the

²⁰https://upload.wikimedia.org/wikipedia/commons/thumb/3/33/Gamma_abs_3D.png/400px-Gamma_abs_3D.png

trace has zero inner product as follows

$$\begin{aligned}
 x &= \frac{1}{1 + \left(\frac{\omega}{\omega_1}\right)^2}, y = \frac{-\left(\frac{\omega}{\omega_1}\right)}{1 + \left(\frac{\omega}{\omega_1}\right)^2} \\
 &\Rightarrow \left(\frac{1}{1 + \left(\frac{\omega}{\omega_1}\right)^2} - 0, \frac{-\left(\frac{\omega}{\omega_1}\right)}{1 + \left(\frac{\omega}{\omega_1}\right)^2} - 0 \right) \cdot \left(\frac{1}{1 + \left(\frac{\omega}{\omega_1}\right)^2} - 1, \frac{-\left(\frac{\omega}{\omega_1}\right)}{1 + \left(\frac{\omega}{\omega_1}\right)^2} - 0 \right) \\
 &= \frac{1}{1 + \left(\frac{\omega}{\omega_1}\right)^2} \left(\frac{1}{1 + \left(\frac{\omega}{\omega_1}\right)^2} - 1 \right) + \frac{-\left(\frac{\omega}{\omega_1}\right)^2}{\left(1 + \left(\frac{\omega}{\omega_1}\right)^2\right)^2} \\
 &= \frac{1}{1 + \left(\frac{\omega}{\omega_1}\right)^2} \frac{1}{1 + \left(\frac{\omega}{\omega_1}\right)^2} - \frac{1 + \left(\frac{\omega}{\omega_1}\right)^2}{\left(1 + \left(\frac{\omega}{\omega_1}\right)^2\right)^2} + \frac{\left(\frac{\omega}{\omega_1}\right)^2}{\left(1 + \left(\frac{\omega}{\omega_1}\right)^2\right)^2} = 0
 \end{aligned} \tag{48}$$

7.2.2. Trace Example: Proportional-Integral (PI) Regulator

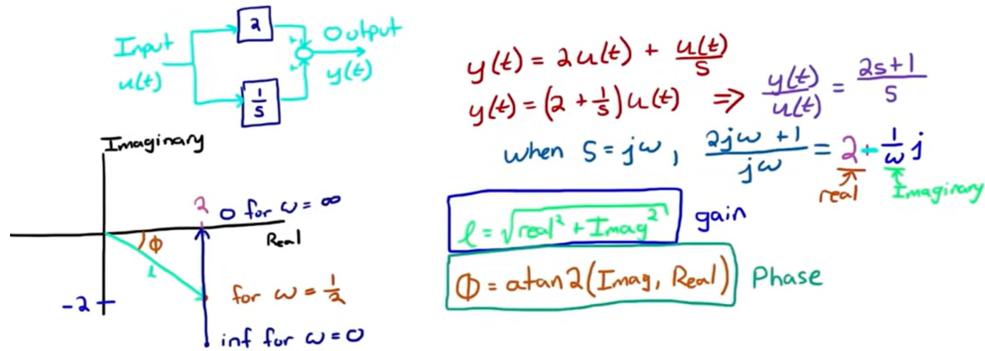


Figure 28. Screenshot from Youtube video: “Control System Lectures - Bode Plots, Introduction”.

Consider the following transfer function suggested by Douglas:

$$C(s) = 1/s + 2$$

which consists of an ideal integrator and a proportional gain. The analysis of the amplitude gain and phase shift and the transfer function trace are shown in Fig. 28.

7.2.3. Trace Example: Proportional Regulator

A constant gain has a transfer function as follows

$$C(s) = K$$

which is a single point in Bode plot, but the sign of K will affect the phase shift.

7.3. Bode Plot of Typical Systems

Drawing the complex values of a transfer function in a s -plane with varying frequency does not explicitly show the frequency. It is desired to draw amplitude gain and phase shift as function of frequency.

Bode plot visualizes the relation between steady state frequency response $T(j\omega)$ versus frequency. Since $T(j\omega)$ is a complex number, two separate plots are needed to show logarithmic gain and the phase shift versus the frequency ω . The logarithmic gain is measured in decibel or dB, which is defined as the common logarithm of the frequency response gain squared:

$$1 \text{ dB} = 10 \log_{10} |T(s)|^2 = 20 \log_{10} |T(s)| \quad (49)$$

If you go look up decibel in Wikipedia, it is defined as a relative value. In this course, this relative value definition only makes sense when decibel is used to describe the **bandwidth** of a system, in which the gain at low frequency for common closed loop control system is $T(0) = 0$ dB.

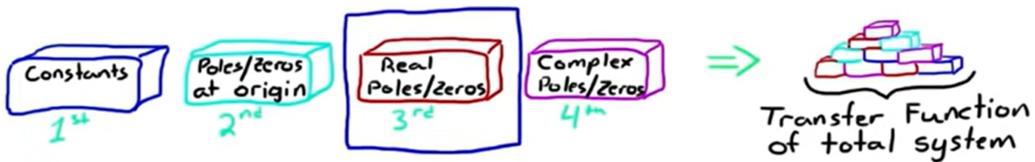


Figure 29. Screenshot from Youtube video: “Bode Plots by Hand: Real Poles or Zeros”.

There are four basic building blocks for transfer functions, as shown in Fig. 29. We are now going to derive the exact curves for each of the building blocks.

7.3.1. Complex Poles/Zeros

Complex poles always appear in the form of conjugate pairs. This is because the model of a physical system has characteristic equation $\text{Den}(s)$ with real coefficients.

Consider the following standard second order transfer function:

$$T(s) = \frac{\omega_0^2}{s^2 + 2\zeta\omega_0 s + \omega_0^2} = \frac{1}{\frac{s^2}{\omega_0^2} + 2\zeta\frac{s}{\omega_0} + 1} \quad (50)$$

where ω_0 is the natural frequency, and note the numerator ω_0^2 was put there to have a dc gain of $T(0) = 1$. Consider steady state, i.e., $s = j\omega$, and treat $T(j\omega)$ as a complex

number:

$$\begin{aligned}
 T(j\omega) &= \frac{1}{\left[\left(\frac{j\omega}{\omega_0}\right)^2 + 1\right] + j\left(2\zeta\frac{\omega}{\omega_0}\right)} \\
 &= \frac{1}{\left[-\left(\frac{\omega}{\omega_0}\right)^2 + 1\right] + j\left(2\zeta\frac{\omega}{\omega_0}\right)} \frac{\left[-\left(\frac{\omega}{\omega_0}\right)^2 + 1\right] - j\left(2\zeta\frac{\omega}{\omega_0}\right)}{\left[-\left(\frac{\omega}{\omega_0}\right)^2 + 1\right] - j\left(2\zeta\frac{\omega}{\omega_0}\right)} \\
 &\Rightarrow \begin{cases} |T(j\omega)| [\text{dB}] = 20 \log_{10} |T(j\omega)| = -20 \log_{10} \sqrt{\left[\left(\frac{j\omega}{\omega_0}\right)^2 + 1\right]^2 + \left(2\zeta\frac{\omega}{\omega_0}\right)^2} \\ \angle T(j\omega) = \arctan 2 \left(-2\zeta\frac{\omega}{\omega_0}, -\left(\frac{\omega}{\omega_0}\right)^2 + 1\right) \end{cases}
 \end{aligned}$$

Evaluate the above expression for three special cases.

- In case of $\omega \ll \omega_0$, we have

$$\frac{\omega}{\omega_0} \ll 1 \Rightarrow \begin{cases} |T(j\omega)| = -20 \log_{10} 1 = 0 \text{ [dB]} \\ \angle T(j\omega) = \arctan 2 \left(-2\zeta\frac{\omega}{\omega_0}, 1\right) = 0^\circ \end{cases}$$

- In case of $\omega = \omega_0$, it gives

$$\frac{\omega}{\omega_0} = 1 \Rightarrow \begin{cases} |T(j\omega)| = -\frac{20}{2} \log_{10} 4\zeta^2 \begin{cases} < 0 \text{ [dB]}, & 4\zeta^2 > 1 \\ = 0 \text{ [dB]}, & \zeta = 0.5 \\ > 0 \text{ [dB]}, & 4\zeta^2 < 1 \end{cases} \\ \angle T(j\omega) = \arctan 2(-2\zeta, 0) = -90^\circ \end{cases}$$

Note only when $\zeta < 0.5$, the amplitude at the natural frequency ω_0 is larger than 0 dB.

- In case of $\omega \gg \omega_0$, it yields

$$\frac{\omega}{\omega_0} \gg 1 \Rightarrow \begin{cases} |T(j\omega)| = -\frac{20}{2} \log_{10} \left[\left(\frac{j\omega}{\omega_0}\right)^4 + \left(2\zeta\frac{\omega}{\omega_0}\right)^2 \right] = -40 \log_{10} \left(\frac{\omega}{\omega_0}\right) \text{ [dB]} \\ \angle T(j\omega) = \arctan 2 \left(-2\zeta\frac{\omega}{\omega_0}, -\left(\frac{\omega}{\omega_0}\right)^2\right) = -180^\circ \end{cases}$$

In the logarithmic diagram, the slope of the curve is measured as -40 dB/dec , meaning whenever the excited frequency of the sinusoid increase by a factor of 10, the gain decreases -40 dB .

Looking at FIGURE 8.10, we have the following observations.

- The resonant frequency ω_r (at which the amplitude gain peaking occurs) is approaching the natural frequency ω_0 as $\zeta \rightarrow 0$. This can be proved by putting the derivative of $|T(j\omega)|$ w.r.t. ω to zero, leading to

$$\omega_r = \omega_0 \sqrt{1 - 2\zeta^2}, \quad \zeta < \sqrt{(2)/2}$$

which means the peak amplitude occurs when the excitation frequency is equal to the resonant frequency ω_r .

- The slope of the magnitude gain curve is -40 dB/dec when $\omega \gg \omega_0$.
- The slope of the phase shift curve is $45^\circ/\text{dec}$ only when damping ratio ζ is close to 1.

The conjugate pair of zeros is the reciprocal of conjugate pair of poles. There is no need to repeat the above derivation.

7.3.2. Real Pole/Zero

Consider the following transfer function $T(s)$ and its frequency response complex gain $T(j\omega)$:

$$\begin{aligned} T(s) &= \frac{1}{\frac{s}{\omega_0} + 1} \xrightarrow{s=j\omega} T(j\omega) = \frac{1}{\frac{j\omega}{\omega_0} + 1} \\ &= \frac{1}{\frac{j\omega}{\omega_0} + 1} \frac{-\frac{j\omega}{\omega_0} + 1}{-\frac{j\omega}{\omega_0} + 1} \\ &\Rightarrow \begin{cases} |T(j\omega)| [\text{dB}] = 20 \log_{10} |T(j\omega)| = -20 \log_{10} \sqrt{1^2 + \left(\frac{\omega}{\omega_0}\right)^2} \\ \angle T(j\omega) = \arctan 2 \left(-\frac{\omega}{\omega_0}, 1\right) \end{cases} \end{aligned}$$

which can as well be evaluated for three special cases depending on ω as compared to ω_0 .

The results of a real zero can be obtained by taking an inverse of the results of a real pole.

7.3.3. Pole/Zero at Origin

The transfer function of a marginally stable real pole leads to the following derivation

$$\begin{aligned} T(s) &= \frac{1}{s} \xrightarrow{s=j\omega} T(j\omega) = \frac{1}{j\omega} = -j\omega \\ &\Rightarrow \begin{cases} |T(j\omega)| [\text{dB}] = 20 \log_{10} |T(j\omega)| = -20 \log_{10} \sqrt{\omega^2} \\ \angle T(j\omega) = \arctan 2(-\omega, 1) \end{cases} \end{aligned}$$

7.3.4. Constant

The results of a constant gain are:

$$T(s) = K \Rightarrow \begin{cases} |T(j\omega)| [\text{dB}] = 20 \log_{10} |T(j\omega)| = -20 \log_{10} K \\ \angle T(j\omega) = \arctan 2(0, K) \end{cases}$$

7.3.5. Asymptotic Curves for Sketching Bode Plot By Hand

See [1, Table 8.1] for list of asymptotic curves for basic building blocks for complicated Bode plot. If a complicated transfer function consists of only four basic building blocks, it is then possible to draw Bode plot by hand by carefully applying slope contribution of each pole or zero from low frequency to high frequency.

Note when the poles and the zeros are very close to each other, the asymptotes approach is no longer valid. Some correction must be made.

7.4. Frequency Response Measurement

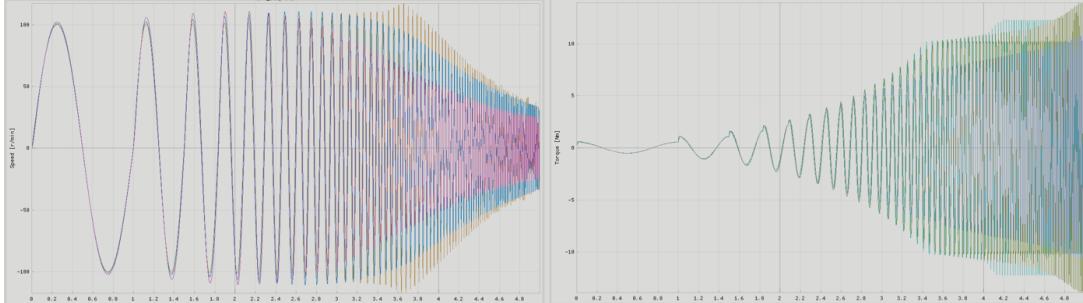


Figure 30. Simulation of a sweep frequency test of a closed loop speed controlled ac motor drive with three sets of tunings. Left: speed response. Right: motor torque.

Sinusoidal is a unique input signal that keeps its shape while passing through an LTI system at steady state with differences from the input only in amplitude and phase angle. This fact has made sinusoid a good test signal for system identification and closed loop control performance measurement.

We can draw Bode plot by hand given the transfer function. In turn, it is also possible for us to deduce the transfer function from the measured frequency response of a system. There are signal analyzer instruments working on this principle. An example is provided in [1, Section 8.3].

Frequency response can also be used to measure the closed loop control performance measurement. For example, a series of sinusoids of different frequencies can be applied to the closed loop system, which is known as the sweep frequency test. Fig. 30 shows the sweep frequency test of the same control system with three different tunings, and the shape of the envelope of the speed response indicates the amplitude gain.

7.5. Performance Specifications in Bode Plot

Performance specification for a standard second order system is specified in terms of rise time and overshoot that are related to natural frequency ω_n and damping ratio ζ . We can have an estimate of ω_n and ζ from the Bode plot.

Natural frequency ω_n is related to the -3 dB bandwidth ω_B . **Bandwidth** ω_B is defined as the frequency at which the frequency response amplitude gain $|T(j\omega)|$ has declined 3 dB (about $1/\sqrt{2}$) from its low-frequency value $|T(0)|$. There also exists definition of bandwidth in terms of phase shift.²¹

Damping ratio ζ is related to amplitude peak value $M_{p\omega}$:

$$M_{p\omega} = |T(j\omega_r)| = \frac{1}{2\zeta\sqrt{1-\zeta^2}}, \quad \zeta < 1/\sqrt{2} \quad (51)$$

²¹Note the definition of bandwidth as drop from peak by 3 dB in signal processing that associates with band-pass filter is completely different from the bandwidth definition in control theory, see <https://electronics.stackexchange.com/questions/280425/bandwidth-of-a-system> for a related discussion. In short, they are completely two things. A simple counter example is a standard second order system with damping ratio less than 0.5 .

Finally, steady state error can also be evaluated by either looking at the low frequency gain $T(0)$ (if it is a finite value) or the slope of the amplitude versus frequency curve at low frequency.

7.6. Relative Stability

Recall the block diagram reduction formula for a negative unity feedback loop is

$$T(s) = \frac{CP}{1+CP} = \frac{L(s)}{1+L(s)}$$

and note for a positive unity feedback loop it becomes

$$T(s) = \frac{CP}{1-CP} = \frac{L(s)}{1-L(s)}$$

where the loop gain $L(s)$ has been defined as the gain of the loop. We will assume negative feedback in the sequel.

Previously, the relative stability has been evaluated as the distance of the pole away from the imaginary axis using RH criterion. In frequency response analysis, we never really solve for the real parts of poles, so we need to develop a different approach to evaluate the relative stability.

7.6.1. Motivation

It is possible to determine the stability of $T(s)$ by evaluating its frequency response complex gain $T(j\omega)$. It is further possible to evaluate the stability of a closed loop complex gain $T(j\omega)$ by looking at its open loop complex gain $L(j\omega)$. If open loop transfer function gain $L(j\omega) = C(j\omega)P(j\omega) \neq -1, \forall \omega > 0$, then the negative feedback control system shall not experience **undamped oscillation** (i.e., **not BIBO stable**).

Note in the complex plane, $-1 = 1\angle-180^\circ$ corresponds to a point $(-1, 0j)$, and $L(j\omega)$ is nothing more than a complex number. If we eliminate the frequency axis of the Bode plot of $L(j\omega)$, then we can obtain a Lissajour-curve-like plot for the frequency response, as exemplified in Fig. 31. We can called it log-magnitude-versus-phase plots, and typical plots for simple transfer functions are summarized in Fig. 32.

7.6.2. Phase Margin and Gain Margin

From Fig. 31, the stability margin (of frequency response) has been defined as the closest distance from the $L(j\omega)$ -trace to the point $(-1, 0j)$ or $(0 \text{ dB}, 0^\circ)$ in the magnitude gain–phase shift plot.

The gain margin indicates how much gain the open loop system can have before the closed loop system becomes not stable. The phase margin shows how much phase shift or time delay the open loop control system can have before the closed loop system becomes not stable.

The Bode plot can also be used to determine if a system is stable by checking if its gain margin and phase margin in two separate plots are positive, as shown in Fig. 33. The gain margin is easily understood: if we add another constant gain of $K = 15 \text{ dB}$ to L_1 in Fig. 33, then the closed loop control system becomes not stable. The phase margin is less intuitive, but as more constant gain is added to L_1 , the phase

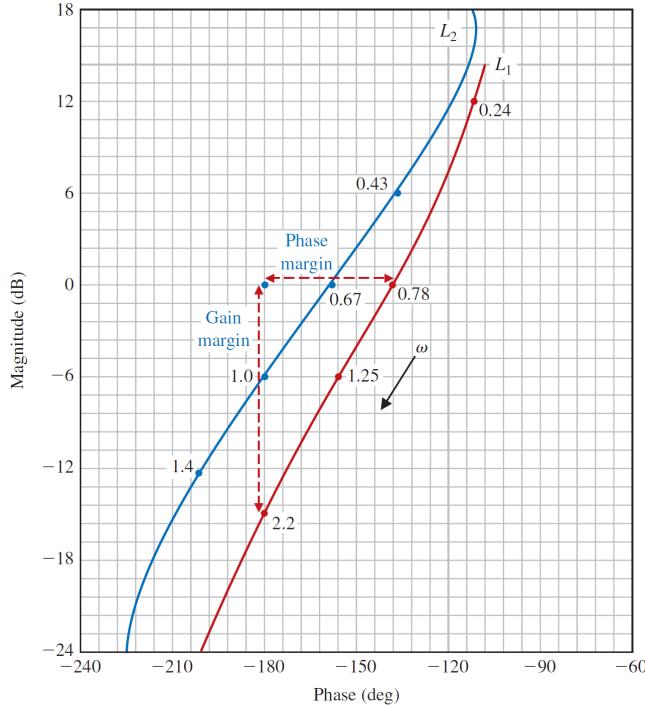


Figure 31. The “Lissajous plot” of the frequency responses of two example systems [1]. The relative stability can be compared in terms of stability margins and L_1 is more stable than L_2 .

margin is found to be less, this is simply because of the monotonically decreasing of the phase-frequency curve.

Stability margin helps to stabilize the system when the actual system has high uncertainty including parameter change or unmodelled dynamics. Previously, we have used the RH criterion for evaluating the relative stability of the system against a vertical line (e.g., $\sigma = -1$) in s -plane. The advantage of a Bode plot over RH criterion is at its ability to visualize relative stability in frequency domain without any iteration for finding the maximal $|\sigma|$.

The following paragraph is quoted from [5, p.467]:

“For satisfactory performance, the phase margin should be between 30° and 60° , and the gain margin should be greater than 6 dB. . . . For minimum-phase systems, the magnitude and phase characteristics of the open loop transfer function are definitely related. The requirement that the phase margin be between 30° and 60° means that in a Bode diagram the slope of the log-magnitude curve at the gain crossover frequency should be more gradual than -40 dB/dec. In most practical cases, a slope of -20 dB/dec is desirable at the gain crossover frequency for stability.”

The design guideline of crossing 0 dB with a slope of -20 dB/dec is in fact very useful and we will have a design example utilizing this guideline. But, what are minimum-phase systems?

7.6.3. Unstable Pole in Open Loop Transfer Function

Unfortunately, this “Bode’s stability criterion” in terms of stability margins only works if the open loop system does not have any zero or pole on the right hand s -plane.

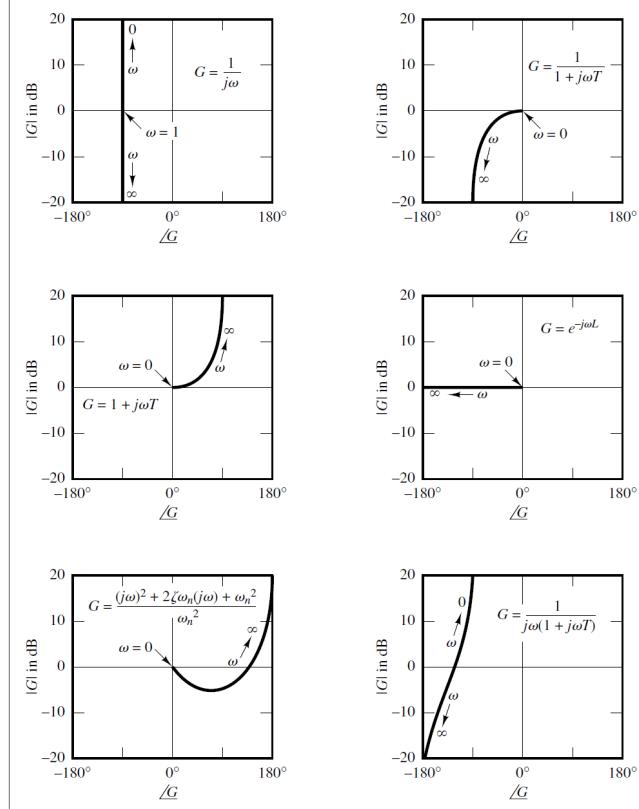


Figure 32. Log-magnitude-versus-phase plots of simple transfer functions [5].

An real life example of an open loop system having RHP pole is seen in the control of power electronic devices to provide power to a constant power load in a dc micro-grid, where the impedance of the load in a local linearization is evaluated to be negative, which results in unstable poles in small signal model.²²

```

1 s = zpk(0, [], 1);
2 unstableT = 100/(s^2 - 4*s + 15);
```

The following MATLAB snippet can be used for checking the stability margins.

```

1 s = tf('s')
2 Popen = 5*(s+3)/s/(s-1)
3 figure
4 allmargin(Popen)
5 bode(Popen)
6 h = bodeplot(Popen);
7 h.showCharacteristic('AllStabilityMargins')
8 P = Popen/(1+Popen)
```

This script show that the open loop system has a negative gain margin ($GM = -14$ dB), but the closed loop system is found to be stable. Therefore, we realize that the gain margin does not have to be positive for a closed loop control system to be stable.

²²See EQUATION (4) and (5) in Emadi et al., Constant Power Loads and Negative Impedance Instability in Automotive Systems: Definition, Modeling, Stability, and Control of Power Electronic Converters and Motor Drives, TTVT VOL. 55, NO. 4, JULY 2006.

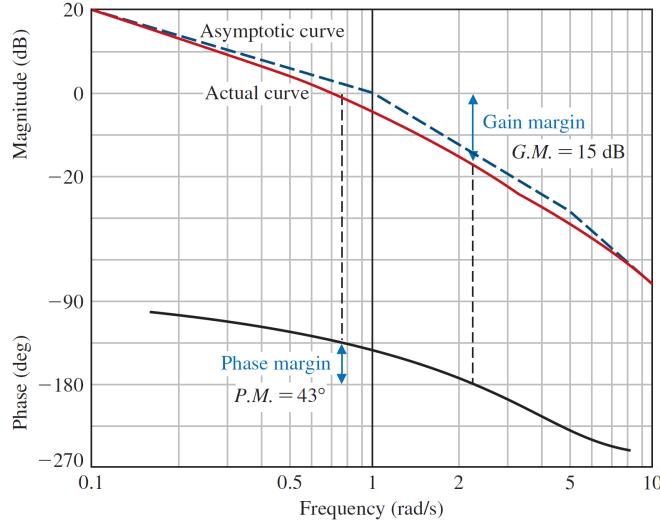


Figure 33. The Bode plot of the frequency responses of the open loop system L_1 from Fig. 31 [1].

For an open loop transfer function that does not have RHP zero/pole, a negative gain margin means its closed loop system is unstable. Therefore, we will need a better tool (Nyquist plot) for evaluating closed loop stability using open loop transfer function.

7.6.4. Relative Stability for Closed Loop System

Recall it has been warned in the beginning of this chapter, that steady state response in frequency domain is not always a correct indicator of stability.

The bode plot of the system having repeated conjugate poles on imaginary axis can be evaluated using MATLAB snippets as follows.

```

1 s = zpk(0, [], 1);
2 P = 1/(s^2+16)^2
3 allmargin(P)
4 h = bodeplot(P);
5 h.showCharacteristic('AllStabilityMargins')

```

Looking at the bode plot of this unstable system, it is a bit difficult to determine if this system is unstable by looking at its stability margins. In fact, the stability margin matters only if we are including this system into a closed loop feedback control system. In fact, the only useful information is the break/corner frequency and the phase shift occurs at that frequency. The sudden change in phase is 360 degrees, implying four poles at the corner frequency.

The bode plot can be plotted for both open loop transfer function and closed loop transfer function. But, the phase margin and gain margin are only valid for analyzing the stability of closed loop system in terms of open loop transfer function. If the relative stability of a closed loop transfer function is of interest, one should resort to the RH criterion that checks how far the poles are away from the imaginary axis.

To recap,

- the distance between the closed loop pole location and imaginary axis is a measure of relative stability of closed loop system;

- the stability margins (i.e., gain margin and phase margin) of the open loop transfer function is also a measure of relative stability of closed loop system;
- the stability margins (i.e., gain margin and phase margin) of the closed loop transfer function is generally not as useful.

7.6.5. RHP Zero in Open Loop Transfer Function

Like unstable pole, a RHP zero will also make “Bode’s stability criterion” fail.

Consider the following open loop system having a RHP zero.

```

1 close all; cla; clc
2 s = zpk(0, [], 1);
3 Popen = 500*(s-2)/(s+1)/(s^2+30*s+229)
4 P = Popen
5 %P = Popen/(1+Popen)
6 %P = Popen/(1-Popen)
7 figure
8 allmargin(P)
9 h = bodeplot(P);
10 h.showCharacteristic('AllStabilityMargins')
11 figure
12 step(P)

```

The step response of the open loop system is stable. However, the closed loop system is unstable, neither via negative feedback nor positive feedback. Therefore, we will need a better tool (Nyquist plot) for evaluating closed loop stability using open loop transfer function.

In addition, for an open loop RHP zero to become an unstable pole in closed loop system $T(s) = CP/(1 + CPH)$, the only way is to have RHP zero in the feedback channel, e.g., $H(s) = \frac{s-1}{s+10}$. This is seldom seen, though.

7.7. Vector Interpretation of the Frequency Response Complex Gain

To understand what is so special about RHP zero and pole, we are now going to introduce a vector tool in pole-zero map.

The contribution by each pole or zero to the frequency response complex gain is not easily examined in the Bode plot in which effects **all** poles and zeros are taken into account for frequency of interest. If we view frequency response of the transfer function as a complex gain, each pole/zero explicit term represents a factor of complex gain that is related to the vector lengths and angles between the pole/zero location and the “responsive point” $(0, j\omega)$, as exemplified in Fig. 34.

If there is a zero on the imaginary axis, it is possible to prevent a certain frequency from passing through the system. An example system is [1, EXAMPLE 8.4]

$$T(s) = \frac{\left(\frac{s}{\omega_0}\right)^2 + 1}{\left(\frac{s}{\omega_0}\right)^2 + 4\frac{s}{\omega_0} + 1}$$

Its gain at $\omega = \omega_0$ is zero, which is contributed by the pair of the conjugate zero.

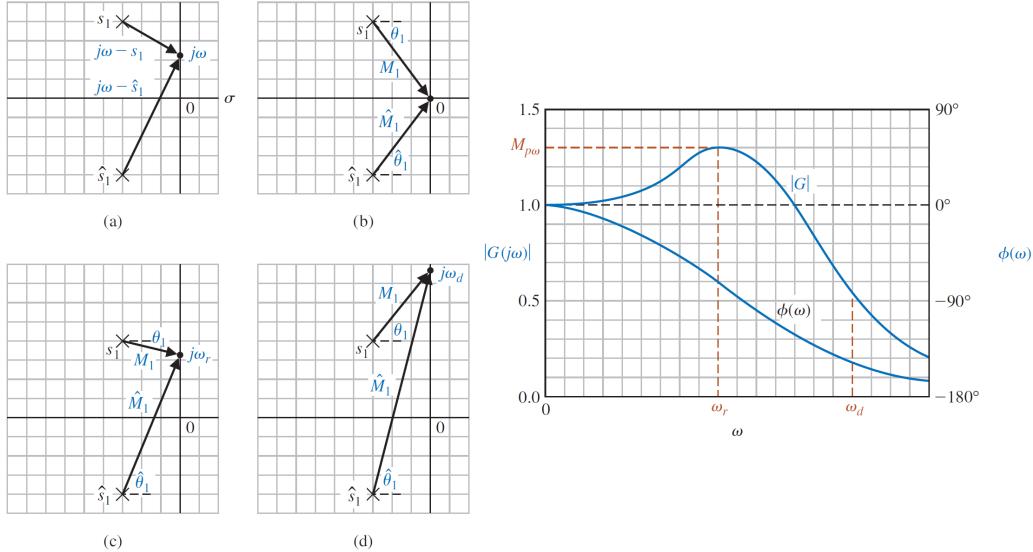


Figure 34. Vector evaluation of a transfer functions in terms of its poles [1].

7.8. Nonminimum Phase System

The vector interpretation of each factor of the complex gain $T(j\omega)$ is found very convenient for analyzing the phase shift contributions of the right hand s -plane zeros. See FIGURE 8.16 in [1].

If we compare systems with right hand s -plane zeros or poles to the one that does not, the Bode plot will report a nonminimum phase shift. Hence, the name. This can be geometrically interpreted because the phase shift is read from the real axis to the frequency responsive point $(0, j\omega)$ in a CCW direction.

An example of realistic nonminimum phase system is the Buck converter with current source supply.²³ The most remarkable indicator of a nonminimum phase system is its step response will have an opposite response against the command at first and then it follows the command.

The definition of nonminimum phase system in the textbook [1] is incomplete. The correct definition is found in other materials, e.g., [5]: “If all the poles and zeros of a system lie in the left half s -plane, then the system is called minimum phase. If a system has at least one pole or zero in the right-half s -plane, then the system is called nonminimum phase.”

7.9. Discussions

Let's discuss applications of Bode plot.

7.9.1. Design Example: Tuning of A Nested Loop Control System

This is an example design from TI's InstaSPIN motion control user guide.²⁴ Consider the nested loop controller for field oriented controlled ac motor shown in Fig. 35.

²³Li et al., *ON Effect of Right-Half-Plane Zero Present in Buck Converters With Input Current Source in Wireless Power Receiver Systems*. TPEL VOL. 36, NO. 6, JUNE 2021.

²⁴www.ti.com/lit/ug/spruhj1h/spruhj1h.pdf

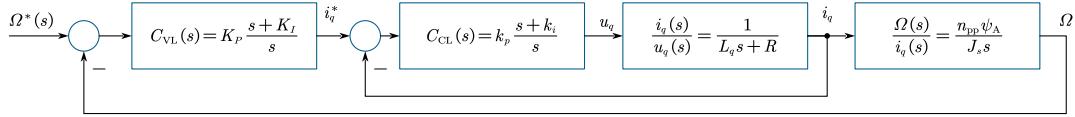


Figure 35. The nested loop controller for field oriented control of ac motor. The Park transform is omitted.

Why nested loop control? If you want to control the speed of the motor, you essentially need to control the time derivative of the speed, and modify the derivative of the speed through the variables that appear in it that you can control. Therefore, we need to take the derivative of the speed, and after the calculation, we find that current appears. We generally do not have a current source power supply, so we need to take the derivative of the time derivative of the speed again, and we see that voltage appears. So we say that the relative order from voltage to current is 2.

- Generally, we need to measure the speed of the motor and calculate the given voltage through the PD control law. D is added to adjust the damping of the second-order system.
- Nested loop control, on the other hand, can convert a control problem of relative order 2 into two control problems of relative order 1, at the cost of adding additional sensors to measure the intermediate state, i.e., the current (which has relative order of 1 from input voltage to the output speed).

Let's apply the loop reduction technique to derive i_q/i_q^* . The inner current loop turns out to be a second order transfer function. The extra pole is introduced by the introduction of the integral control term $k_p k_i / s$. In practice, the inner loop can take advantage of this extra zero, and impose a **zero-pole cancellation**. To this end, we need to place the zero ($s = -k_i$) of the open loop transfer function to R/L_q [rad/s]. After zero-pole cancellation, we define the corner frequency of the rest first-order transfer function

$$\frac{i_q}{i_q^*} = \frac{1}{\frac{s}{CLBW} + 1}$$

as the -3 dB bandwidth of current closed loop system, i.e.,

$$k_p = L_q \times CLBW$$

where CLBW [rad/s] stands for current loop bandwidth. This is not a standard second order transfer function because of the presence of the zero at $s = -K_I$ introduced by the integral control. We cannot use dominant poles or use transient performance metrics.

The speed closed loop transfer function is

$$\frac{\Omega(s)}{\Omega^*(s)} = \frac{\frac{K_P n_{pp} \psi_A}{J_s s^2} \frac{s + K_I}{\frac{s}{CLBW} + 1}}{1 + \frac{K_P n_{pp} \psi_A}{J_s s^2} \frac{s + K_I}{\frac{s}{CLBW} + 1}} = K_P \frac{n_{pp} \psi_A}{J_s} \frac{s + K_I}{\frac{1}{CLBW} s^3 + s^2 + K_P \frac{n_{pp} \psi_A}{J_s} (s + K_I)} \quad (52)$$

The following design with another tuning button called damping factor δ [1]:

$$\begin{aligned} k_i &= \frac{R}{L_q} \\ k_p &= L_q \times \text{CLBW} \\ K_I &= \frac{\text{CLBW}}{\delta^2} \\ K_P \frac{n_{pp}\psi_A}{J_s} &= \delta K_I = \frac{\text{CLBW}}{\delta} \end{aligned} \quad (53)$$

is suggested to make sure the amplitude gain of the speed open loop frequency response passes through 0 dB at the peaking of phase shift- ω curve, such that the largest phase margin can be approximately obtained. This design works because the presence of the zero between two poles. The first pole is

The following snippet shows a design that does not fully follow the above K_P and K_I tuning rule so the phase margin is not as large as possible.

```

1 close all; cla; clc
2 s = zpk(0, [], 1);
3 Popen = 100/s^2 * (s+10) / (s/100+1)
4 P = Popen/(1+Popen)
5 allmargin(Popen)
6 h = bodeplot(Popen);
7 h.showCharacteristic('AllStabilityMargins')
8 figure; bode(P)
9 figure; step(P)

```

As a comparison, the following snippets implement the above tuning rule (53). The results are shown in Fig. 36

```

1 close all; cla; clc; s = zpk(0, [], 1);
2 % two tuning buttons
3 CLBW = 200; delta = 2;
4 % inner loop
5 Lq = 5e-3; R = 1; kp = CLBW*Lq; ki = R/Lq;
6 PCL = 1/(Lq*s+R); CCL = kp*(1+ki/s);
7 PclosedInner = CCL*PCL / (1+CCL*PCL)
8 % outer loop
9 n_pp = 4; psi_A = 0.1; Js = 0.006; K = n_pp * psi_A / Js;
10 KP = CLBW / delta / K; KI = CLBW/delta^2;
11 PVL = K/s; CVL = KP*(1+KI/s);
12 Popen = CVL * PclosedInner * PVL;
13 PopenMargin = allmargin(Popen)
14 P = Popen/(1+Popen)
15 P = minreal(P)
16 subplot(221); h = bodeplot(Popen); h.showCharacteristic(
    'AllStabilityMargins')
17 subplot(222); bode(P); grid; h1 = findobj(gcf,'type','line'); set(
    h1,'linewidth',2);
18 subplot(223); step(P); grid; h1 = findobj(gcf,'type','line'); set(
    h1,'linewidth',2);
19 subplot(224); pzmap(P); grid; h1 = findobj(gcf,'type','line'); set(
    h1,'linewidth',2);

```

The closed loop speed control system has three poles and one zero, and as a result, its Bode plot in Fig. 36 is very similar to a second order system. Generally speaking,

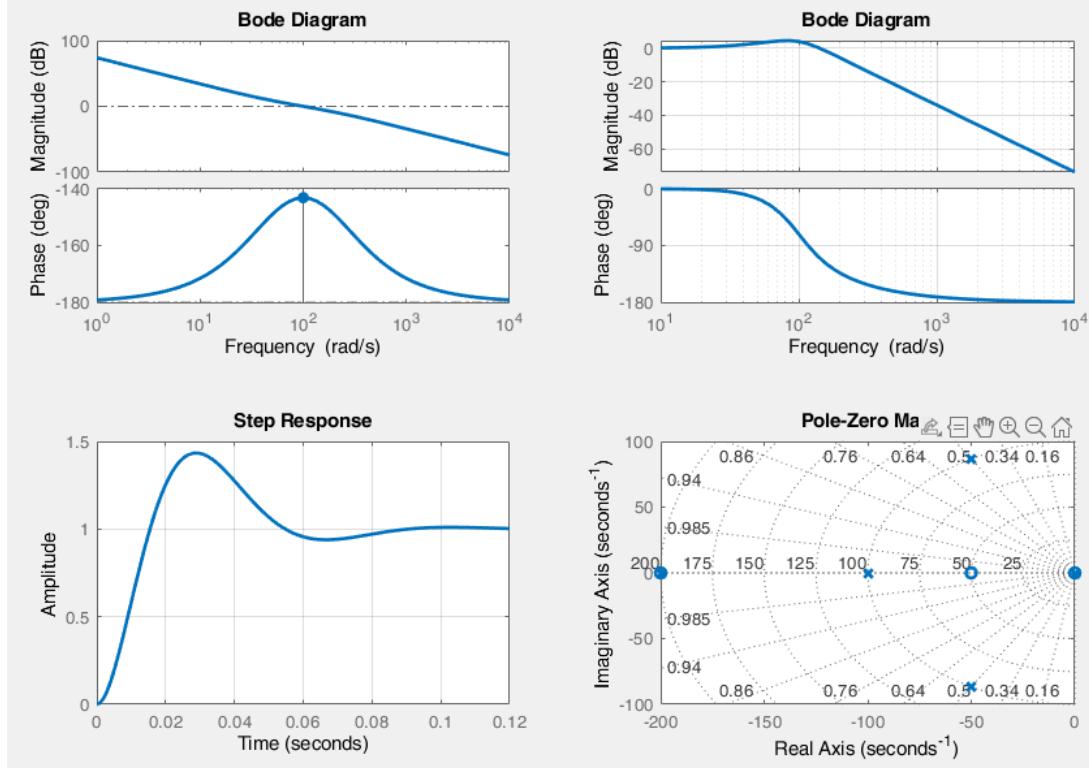


Figure 36. Nested loop control design with CLBW = 200 rad/s and $\delta = 2$. For motor parameters, see code snippets.

a zero is able to cancel the effect of a pole, as long as they are close enough in the complex s -plane.

The transfer function from load torque disturbance $-T_L(s)$ to speed output $\Omega(s)$ can be analyzed using the following code snippet.

```

1 openloopDisturbance = 1/s / Js;
2 H = CVL * PclosedInner;
3 PDisturbance = openloopDisturbance/(1+openloopDisturbance*H)
4 figure; bode(PDisturbance); grid; h1 = findobj(gcf,'type','line');
    set(h1,'linewidth',2);

```

7.9.2. Frequency Response of The Gang Members

Recall the gang members from Section 4.8 defined for evaluation of different inputs. The gang members' frequency responses are important for practical control system design.

Let's consider a simple example with the aid of Matlab.

```

1 P = tf([1], [1, 1])
2 C = tf([1, 1], [1 0])
3 subplot(141); bode(C * P/(1+C * P)); grid; h1 = findobj(gcf,'type',
    'line'); set(h1,'linewidth',2);
4 subplot(142); bode(1/(1+C * P)); grid; h1 = findobj(gcf,'type',
    'line'); set(h1,'linewidth',2);
5 subplot(143); bode(P/(1+C * P)); grid; h1 = findobj(gcf,'type',
    'line'); set(h1,'linewidth',2);

```

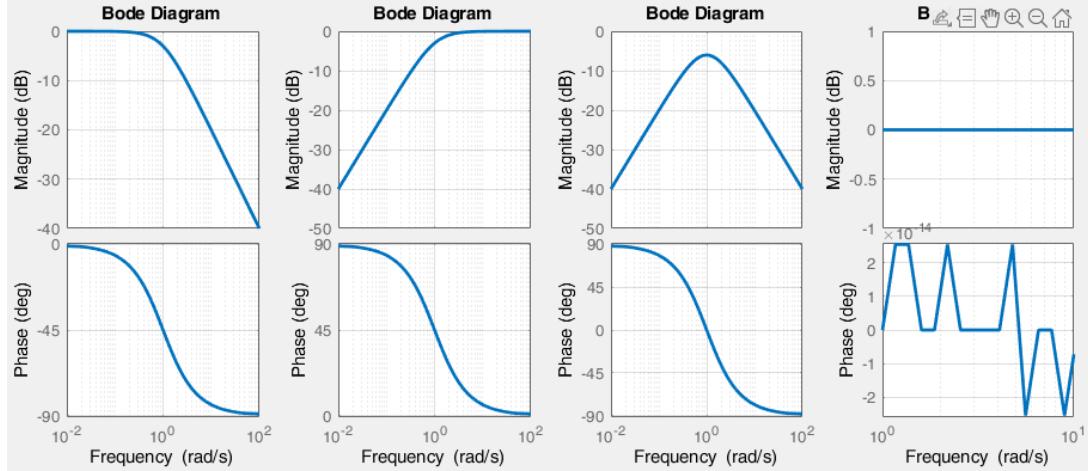


Figure 37. Sensitivity function frequency response. From left to right are noise to error $1 - S$, reference to error PS , and CS (which equals to 1).

```
6 subplot(144); bode(C/(1+C * P)); grid; h1 = findobj(gcf,'type','line'); set(h1,'linewidth',2);
```

where the controller $C(s) = (s + 1)/s$ is a PI regulator.

7.9.3. Bode Plot of the Closed Loop System

When discussing relative stability, it is assumed we want to check the stability of the closed loop system by looking at the Bode plot of the open loop frequency response. The Bode plot can of course be directly used for closed loop system, and it is possible to distinguish between a stable closed loop transfer function and an unstable closed loop transfer function by looking at their Bode plots.

8. Stability in Frequency Domain

We have learned RH criterion applied to characteristic equation and Bode plot of frequency response. A comparison between these two methods shall now be put forth.

Relative stability has distinct meanings in those two methods.

- The relative stability by RH criterion is evaluated by the distance of the roots of the characteristic equation from the $j\omega$ axis. Since only the sign of the real parts of the poles can be evaluated by RH criterion, this distance can only be found through a trial and error approach, which is deemed inconvenient.
- The relative stability by steady state frequency response is evaluated by the horizontal and vertical distance of the complex gain curve from the $(0 \text{ dB}, -180^\circ)$ point in the Lissajour-curve-like plot of the frequency response. These two distances are known as stability margin (i.e., gain margin and phase margin). Stability margins read from Bode plot are only valid for minimum phase open loop system.

Note the key difference is that the Bode plot looks into the open loop system to figure out the relative stability of the closed loop system, while the RH criterion directly checks the characteristic equation of closed loop system to check relative stability.

In terms of sufficiency and necessity, they are different, because the complex gain of the steady state frequency response is not equivalent to the characteristic equation.

- RH criterion based on characteristic equation is a sufficient and necessary condition of BIBO stability, simply because it checks for stability by counting the number of right hand s -plane poles. The sufficient and necessary condition for stability is not having RHP roots in the characteristic equation (of the closed loop transfer function).
- The Bode plot method, on the other hand, does not really tell us information about the unstable roots of the characteristic equation (of a closed loop function). The Bode plot stability criterion is motivated by the fact that the steady state amplitude gain of the closed loop system should not become infinity for sinusoidal input of any frequency.

For this reason, we should stick with the RH criterion for checking stability, but Bode plot is more convenient to check relative stability of a closed loop control system whose open loop transfer function is a minimum phase system that has no RHP zeros/poles.

This chapter is going to introduce a tool called **Nyquist plot**, and it is a result of mapping Nyquist contour to another complex plane. This is the tool that checks RHP poles of the closed loop system by looking at the roots of the characteristic equation, and at the same time, it allows to visually check stability margins.

8.1. Motivation and Prerequisite

We assume the zeros and poles of the plant $P(s)$ and sensor $H(s)$ are known, and of course we are aware how many poles there are in the controller $C(s)$. The bad news is that, the characteristic equation $1 + CPH = 0$ is difficult to solve, especially when there are new poles introduced in $C(s)$ and $H(s)$.

The whole chapter makes sense only you agree that $1 + CPH = 0$ cannot be solved easily and the poles of $C(s)P(s)H(s)$ are already known.

8.2. Nyquist Plot

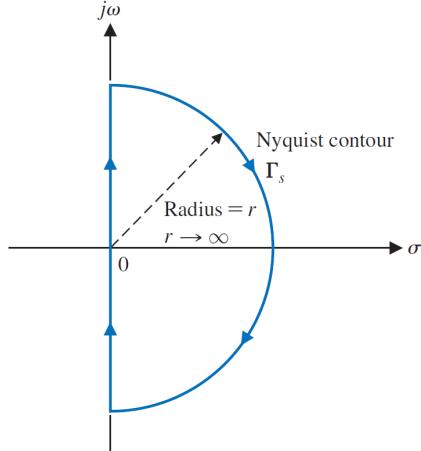


Figure 38. The Nyquist contour Γ_s . The Nyquist contour starts at origin and covers the whole right hand s -plane in a clockwise direction. The infinitely large half circle will be mapped to a point in the image plane, i.e., the $L(s)$ -plane.

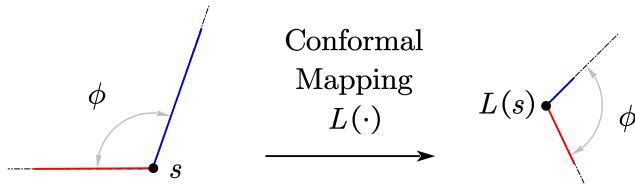


Figure 39. Conformal mapping with a flipped mapping in the local. The angle between the two local line segments that cross at point s is preserved after the conformal mapping.

I came across the lecture notes of ECE 486 from UIUC.²⁵ I suggest you to take a look, and the key take-away's are summarized below.

- (1) Nyquist contour Γ_s in Fig. 38 begins at origin, moves vertically along the imaginary axis to infinity. This happens to be the same as we evaluate the frequency response by sweeping $s = j\omega$ with $\omega \in [0, \infty]$. Therefore, this method is said to evaluate stability in the frequency domain.
- (2) Nyquist plot is a set of points and is the result of mapping Nyquist contour in s -plane to the $L(s)$ -plane, where $L(s)$ is the open loop transfer function $L(s) = CPH$. Nyquist plot of $L(j\omega)$ can be viewed as the image of the imaginary axis $\{j\omega \in \mathbb{C} : -\infty < \omega < \infty\}$ under the mapping $L : \mathbb{C} \mapsto \mathbb{C}$.
- (3) The conformal mapping $L(s)$ is a function of complex number, as long as its derivative with respect to s is not equal to zero. Conformal mapping is angle-preserving. Angle-preserving means that conformal mapping preserves the angle of any two arbitrary smooth curves that cross at point s after the mapping. For example, consider an arbitrary point on a contour, where the left segment to the point and the right segment to the point are two curves, and the angle of these two segment is not changed during a conformal mapping, as shown in Fig. 39. Angle-preserving property ensures that the Nyquist contour preserves its

²⁵Nyquist Stability Criterion in ECE487: <https://courses.engr.illinois.edu/ece486/fa2018/handbook/lec18.html>.

rotating direction in the image plane or $L(s)$ -plane as long as the contour is not flipped after the conformal mapping. To understand how flipping occurs during a mapping, let's consider a point $s = re^{j\phi}$ in the Nyquist contour, the inverse mapping $L(s) = s^{-1}$ flips the contour in the image plane, because $s^{-1} = r^{-1}e^{-j\phi}$ in which a negative sign is seen in the angle/argument of this complex number. On the contrary, the conjugate inverse mapping $L(s) = 1/\bar{s} = r^{-1}e^{j\phi}$ does not flip the contour, i.e., the image contour preserves the rotation direction.

- (4) The argument principle (in control) [1] states that the s -plane closed contour Γ_s (*that does not self-intersect, and does not pass through any poles or zeros*) encircles Z zeros and P poles of $L(s)$ in clockwise direction, and then its image contour $\Gamma_{L(s)}$ will encircle the origin of $L(s)$ -plane $N_{cw} = Z - P$ times in the clockwise direction, and negative N_{cw} indicates the number of encirclement in contour-clockwise direction, i.e., $N_{ccw} = -N_{cw}$.²⁶
- (5) Nyquist stability criterion states that closed loop system is stable if and only if the Nyquist plot of its characteristic equation $1 + L(s)$ encircles the point $(0, 0)$ counterclockwise P times, where P is the number of unstable open loop poles of $L(s)$, i.e.,

$$\begin{aligned} P &= \#(\text{unstable open loop poles}) \\ &= Z - N_{cw} \\ &= \#(\text{unstable roots}) + N_{ccw} \\ &= 0 + N_{ccw} \end{aligned} \tag{54}$$

This is understood by considering only the angles of the following fraction

$$\begin{aligned} \angle L(s) &= \angle \frac{(s - z_1) \cdots (s - z_m)}{(s - p_1) \cdots (s - p_n)} \\ &= \angle \frac{e^{j\psi_1} \cdots e^{j\psi_m}}{e^{j\varphi_1} \cdots e^{j\varphi_m}} \\ &= \sum_{i=1}^m \angle (s - z_i) - \sum_{j=1}^n \angle (s - p_j) \\ &= \sum_{i=1}^m \psi_i - \sum_{j=1}^n \varphi_j \end{aligned}$$

where note that $L(s)$ consists of physically realizable systems in practice, therefore we have the dimension $n \geq m$. When Nyquist contour encircles a zero or a pole, the change in ψ_i or φ_j will be 360° .

- (6) Nyquist criterion is the frequency domain equivalent of Routh-Hurwitz stability criterion. This implies that the characteristic equation is evaluated for closed loop system stability. The roots are not solved. Only the positiveness of the real parts of the roots are checked.
- (7) Nyquist plot is symmetric about the horizontal real axis, because:

$$L(-j\omega) = \overline{L(j\omega)}$$

²⁶In case you are interested in a formal proof of the argument principle in math, please see https://www.youtube.com/watch?v=79-ESkh5_f0.

where the over-line denotes the conjugate of the complex number.

- (8) A not proper transfer function has a Nyquist plot that ends infinitely far away from origin.
- (9) A proper but not strictly proper transfer function has a Nyquist plot that ends at the real axis when $s = \infty$.
- (10) A strictly proper transfer function has a Nyquist plot that ends at origin.
- (11) Relative stability can be read from the Nyquist plot. The Nyquist plot is closely related to the Lissajour-curve-like trajectory of the Bode plot. The main difference is how you should read the gain margin. In dB-degrees plot, it is the difference between the critical point (0 dB, 180°) and the crossing at the horizontal line:

$$GM [dB] = 0 \text{ dB} - |L(j\omega)|_{\omega=\arg[\angle L(j\omega)=180^\circ]} \quad (55)$$

In s -plane, gain margin is read as a division:

$$GM = 1 / |L(j\omega)|_{\omega=\arg[\angle L(j\omega)=180^\circ]} \quad (56)$$

- (12) The proportional gain in the controller $C(s)$ acts as a scale to the Nyquist plot. Therefore, if the Nyquist plot has some portion on the left hand $L(s)$ -plane, it has a finite gain margin GM.
- (13) When there is a pole $L(s)$ on the imaginary axis, the gain becomes infinite and the phase becomes indefinite when $s = j\omega$ passes through the zero. In this case, we need to adjust the Nyquist contour around this open loop marginally stable pole, such that the phase in the $L(s)$ -plane can be determined. Think what happens when there is repeated poles at origin?
- (14) When there is a zero $L(s)$ on the origin, there is no need to adjust the Nyquist contour. The Nyquist plot goes to 0 when Nyquist contour passing the marginally stable open loop zero.

8.2.1. Steps to Plot a Nyquist Plot

Steps to plot the Nyquist plot are now summarized.

- (1) If a numerical solution is needed, sketch the Bode plots by hand to obtain the magnitude versus frequency plot and phase versus frequency plot. Use Bode plot to sketch Nyquist plot. Pay attention to the crossings of 0 dB line and 90°, 180°, 270°, 360° lines. The 0 dB crossing point is mapped to the unit circle; while crossings of 90°, 180°, 270°, 360° are mapped to the real and imaginary axes of the image plane.
- (2) If a symbolic solution is needed, compute the real part and imaginary part of $L(j\omega)$.
 - The crossing point between the Nyquist plot and the real axis can be obtained by first setting $\text{Im}[L(j\omega)] = 0$ to obtain the angular frequency $\omega_{\text{Im}=0}$, and then compute for $\text{Re}[L(j\omega_{\text{Im}=0})]$. The crossing point is $(\text{Re}[L(j\omega_{\text{Im}=0})], 0)$.
 - The crossing point between the Nyquist plot and the imaginary axis can be obtained by first setting $\text{Re}[L(j\omega)] = 0$ to obtain the angular frequency $\omega_{\text{Re}=0}$, and then compute for $\text{Im}[L(j\omega_{\text{Re}=0})]$. The crossing point is $(\text{Re}[0, L(j\omega_{\text{Im}=0})])$.

- (3) If there is a pole at the origin, to be consistent with the “nyqlog” function, we choose to exclude the marginally stable pole from the Nyquist contour. See EXAMPLE 9.2 and EXAMPLE 9.3 in the textbook [1].
- (4) Count the encirclements, count unstable open loop poles, and apply Nyquist criterion $Z = N_{\text{cw}} + P$ to find the unstable roots of the characteristic equation.

The following MATLAB snippet calculates the frequency response of a transfer function at 10 rad/s.

```
1 evalfr(L, 10j)
```

The following MATLAB plots only the portion of a Nyquist plot when $\omega > 0$.

```
1 plotoptions= nyquistoptions('cstprefs') ; % default options
2 plotoptions.ShowFullContour = 'off'; % exclude negative
   frequencies
3 nyquist(L, plotoptions)
```

Alternatively, the MATLAB function nyqlog written by a third party provides indicator for positive frequencies and negative frequencies as well. Moreover, it handles infinite long Nyquist plot well by using a logarithm plot.²⁷

```
1 nyqlog(L)
```

Trick question 1: What happens when the Bode plot does not intercept 180° ? Examples include: $L = 1/(s+1)$, and $L = s/(s+2)/(s-5)$. Hint: the phase does not matter if the magnitude is zero.

Trick question 2: What is the rotating direction when the Nyquist contour encircles a pole and a zero at the same time? Check for $L = (s-2)/(s-5)$, $L = (s-5)/(s-2)$, $L = (s-5)/(s+2)$, and $L = (s-2)/(s-5)/(s-1)$ to find out. Hint: one unstable pole flips the Nyquist plot once.

If it is required to use the Nyquist plot for determining the range of a gain $K \geq 0$, the key step is to calculate for the crossing point of the Nyquist plot to the negative real axis. **Compare the crossing point's coordinate with $(-1, 0)$ to determine the range of K .**

Example 1: find the range of K that makes

$$1 + L(s) = 1 + K/(s+1)/(s+2)$$

stable. Plot a stable Nyquist plot for a value of K . Hint: the point $(-1/K, 0)$ must avoid getting inside the Nyquist plot.

Example 2: find the range of K that makes

$$1 + L(s) = 1 + K / (s-1) / (s^2 + 2s + 3)$$

stable. Plot a stable Nyquist plot for a value of K . Find the gain margin and phase margin.

Example 3: find the range of K that makes

$$1 + L(s) = 1 + K (s-1) / (s+2) / (s^2 - 2s + 1)$$

²⁷Don't download version 1.6 (2016), but download version 1.5 (2009) instead. https://www.mathworks.cn/matlabcentral/fileexchange/7444-nyquist-plot-with-logarithmic-amplitudes#version_history_tab. Or, just go to course material repository: <https://github.com/horychen/ee160>.

stable. Plot a stable Nyquist plot for a value of K . Find the gain margin and phase margin.

The following MATLAB snippet is useful to find answers to the above examples.

```

1 close all; cla; clc; s = zpk(0, [], 1);
2 L = 1/(s+1)/(s+2)
3 L = 3.5/(s-1)/(s^2+2*s+3)
4 L = 1.75*(s-1)/(s+2)/(s^2-s+1)
5 nyqplot(L)
6 subplot(121)
7 h = bodeplot(L);
8 h.showCharacteristic('AllStabilityMargins')
9 grid; h1 = findobj(gcf,'type','line'); set(h1,'linewidth',3);
10 subplot(222)
11 h=nyquistplot(L);
12 h.showCharacteristic('AllStabilityMargins')
13 h1 = findobj(gcf,'type','line'); set(h1,'linewidth',3);
14 daspect([1 1 1])
15 subplot(224)
16 nyqlog(L)

```

Example 4: Stability and fast responsiveness are conflicting. Consider the pitch system of the auto pilot of F16 example by Brian Douglas, where the input is angle of the elevator, and the output is the pitch angle of the aircraft.

```

1 format long
2 L=(3.553e-15*s^4-0.1642*s^3-0.1243*s^2-0.00161*s+9.121e-17)/(s
   ^5+1.825*s^4+2.941*s^3+0.03508*s^2+0.01522*s-1.245e-15)
3 roots([1 1.825 2.941 .03508 .01522 -1.245e-15])
4 nyqplot(L)
5 subplot(121)
6 h = bodeplot(L);
7 h.showCharacteristic('AllStabilityMargins')
8 grid; h1 = findobj(gcf,'type','line'); set(h1,'linewidth',3);
9 subplot(122)
10 h=nyquistplot(L);
11 h.showCharacteristic('AllStabilityMargins')
12 h1 = findobj(gcf,'type','line'); set(h1,'linewidth',3);
13 daspect([1 1 1])

```

Somehow the produced results are a bit off from the results of the Youtube video (tsgOstfoNhk).

8.3. Nichols Chart*

9. Root Locus Method

The Nyquist plot becomes difficult to read when there are poles on the imaginary axis. For example, try to use the MATLAB script from last chapter to draw the Nyquist plot of $L(s) = 10/(s(s + 1)(s^2/4 + 1))$.

We have analyzed system performance in time domain (in terms of metrics including peak time, rise time, settling time, steady state error) and frequency domain (in terms of metrics including phase margin, gain margin, and bandwidth).

The system stability can be analyzed using RH criterion and Nyquist plot that both focus on the characteristic equation of closed loop system. This chapter focuses on s -domain analysis tool, known as the root locus method. It can be used for analysis of stability and performance.

9.1. Motivation: are these two systems equivalent?

I have said in class that I don't feel like the block diagram reduction is actually useful for practical application. Now I am ready to give a counter-example. Consider two systems that have the same closed loop transfer functions, as shown in Fig. 40a and Fig. 40b, but are they really equivalent, given different open loop transfer functions?

Their difference in the open loop transfer functions results in completely different root loci in the s -plane as evaluating systems in Fig. 40c and Fig. 40d. Sample codes are listed below, where minreal is essential to clean up the root loci by removing pole zero cancellation pairs.

```

1 s = zpk(0, [], 1);
2 CP = 1/s/(s+50);
3 H = 1/(s+30);
4 L1 = CP*H
5 L2 = CP / ( 1 + CP*H - CP); L2 = minreal(L2)
6 figure
7 subplot(221); rlocus(L1)
8 subplot(222); pzmap(L1)
9 subplot(223); rlocus(L2)
10 subplot(224); pzmap(L2)

```

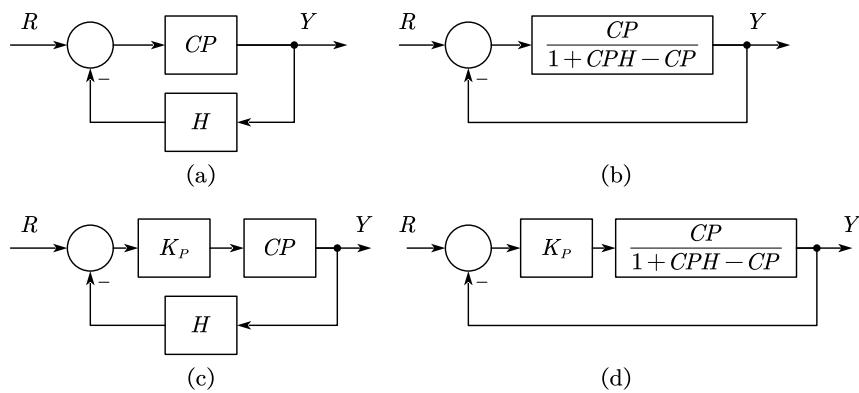


Figure 40. Two systems having the same closed loop transfer function are not equivalent in root locus analysis.

9.2. Tuning Arbitrary Parameter Using Root Locus Method

The root locus method assumes the varying parameter is a constant gain applied to the open loop transfer function. When parameter appears in a random location in the transfer function, can we still use root locus method?

Consider the following example from one of Brian's video²⁸

$$C(s)P(s) = \left(5 + 0.3s + \frac{5}{s\tau}\right) \frac{1}{0.01s^2 + 0.11s + 0.1} \quad (57)$$

The time constant τ for the integral control is the variable to be tuned.

Generally speaking, it is okay to directly calculate for the roots of $1 + CP$ with parameter appearing in arbitrary location, e.g., τ in (57), and plot them in the s -plane, but we lose the intuition about how the roots move in the s -plane.

9.3. Rules to Plot Root Locus

Rearrange the characteristic equation such that the parameter to be tuned appears as a gain denoted by K :

$$1 + L(s) = 1 + K \frac{Q(s)}{P(s)}$$

- (1) Roots move from the poles of $L(s)$ to the zeros of $L(s)$. If there are more poles than zeros, those poles with unmatched “zero friend” will go to infinity to seek their friends. Similarly, if there are more zeros than poles, loci will come from infinity to those zeros with unmatched poles. As a result, there are n loci where n is the degree of $Q(s)$ or $P(s)$ whichever is greater. In most cases, we use n as the order of polynomial $P(s)$, and use M to denote the order of polynomial $Q(s)$, and for physically realizable systems, $n \geq M$.
- (2) Roots that are not on the real axis, always appear in conjugate pairs, because the characteristic equation $1 + L(s)$ has real coefficients.
- (3) Each root locus does not cross over itself.
- (4) The portion of the real axis to the left of an odd number of open loop poles and zeros are part of the loci, because the angles contributions of all zeros and poles add up to 180° when $K \geq 0$.
- (5) A pair of loci leave and enter the real axis at 90° .
- (6) The $n - M$ lines going to infinity have asymptotes that have angles determined by $\phi_A = 180^\circ(2k + 1)/(n - M)$, with $k = 0, 1, 2, \dots, (n - M - 1)$. For example, if there are three lines go to infinity, they go at $60^\circ, 180^\circ$, and 300° degrees; and if there are four lines go to infinity, they go at $45^\circ, 135^\circ, 225^\circ, 315^\circ$ degrees.
- (7) The $n - M$ lines going to infinity have asymptotes centered at the real axis at $(\sigma_A, 0)$. The real scalar σ_A can be calculated by dividing the difference of the sum of poles minus the sum of zeros by $n - M$.
- (8) The gain value K that corresponds to the crossing point of root loci with the imaginary axis can be found by RH criterion.
- (9) The exact departure point from the real axis can be found by using $1 + L(s) = 0$ to define $K = K(s)$ and find the real value of s that makes $K(s)$ reach maximum.

²⁸PID and Root Locus <https://www.youtube.com/watch?v=z1G2sDEG5yQ>

Example: consider the control of the magnetically levitated ball with the plant transfer function as $1/(s^2 - 16)$. Try to explain why adding PD control is helpful for stabilization but adding PI control is not enough. Furthermore, how about adding a PID controller?

9.4. Time Delay and Padé Approximation

In practice, time delay appears as the dead time of actuator. During dead time, the actuator is not responsive regardless of the command. In continuous time domain, dead time can be modelled as time delay.

The time delay in time domain is modelled using delta function $\delta(t - T_d)$. Its Laplace transform is e^{-sT_d} , which can be written as Taylor series with a change of variable $x = sT_d$

$$e^{-x} = \sum_{n=0}^{\infty} \frac{(-x)^n}{n!} = 1 - x + \frac{1}{2}x^2 - \frac{1}{6}x^3 + \dots$$

which is a polynomial having infinite number of states (or zeros if you view e^{-sT_d} as a transfer function).

Root locus cannot handle infinite number of poles or zeros. If you run the following script in MATLAB

```
1 s = zpk(0, [], 1);
2 L = exp(-s) / (s+1) / (s+2)
3 rlocus(L)
```

you will get error message: "Error using DynamicSystem/rlocus (line 65) The "rlocus" command cannot be used for continuous-time models with delays. Use the "pade" command to approximate delays."

In contrast, frequency domain analysis tool works well for time delay, because the time delay can be easily accounted for as the phase delay of frequency response.

The Padé approximation (of $e^{-sT_d} = e^{-x}$) of order [m/n] is the rational function

$$R(x) = \frac{\sum_{j=0}^{\infty} a_j x^j}{1 + \sum_{j=1}^{\infty} b_j x^j} = \frac{a_0 + a_1 x + a_2 x^2 + \dots + a_m x^m}{1 + b_1 x + b_2 x^2 + \dots + b_n x^n}$$

The Padé table summarizes the commonly used approximations, where only the diagonal components introduce pure phase delay, i.e., those have a gain of 0 dB.²⁹

9.5. Revisit Nonminimum Phase System

The definition of nonminimum phase system can be further extended, because there is in fact only one minimum phase system but there are infinite number of nonminimum phase systems.

Consider the following example presented as code snippet.

```
1 %% Delay and nonminimum phase system
2 Td = 1.0; % time delay = 1 sec is fair
3 P = tf([1,2], [1,3,1])
```

²⁹https://en.wikipedia.org/wiki/Pad%C3%A9_table#An_example_%E2%80%93_the_exponential_function

```

4 P_delay = tf([1,2], [1,3,1], 'InputDelay', Td)
5 P_RHP_zero = tf([-1,2], [1,3,1])
6 [n,d] = pade(Td, 2)
7 P_RHP_pade = P * tf(n,d)
8 options = bodeoptions;
9 options.PhaseWrapping = 'on';
10 subplot(221); bode(P, P_delay, P_RHP_zero, P_RHP_pade, options);
    grid; h1 = findobj(gcf,'type','line'); set(h1,'linewidth',2);
11 subplot(222); step(P, P_delay, P_RHP_zero, P_RHP_pade); grid; h1 =
    findobj(gcf,'type','line'); set(h1,'linewidth',2);
12
13 % closed loop response
14 P = P/(1+P)
15 P_delay = P_delay /(1+P_delay)
16 P_RHP_zero = P_RHP_zero /(1+P_RHP_zero )
17 P_RHP_pade = P_RHP_pade /(1+P_RHP_pade )
18 subplot(223); bode(P, P_delay, P_RHP_zero, P_RHP_pade, options);
    grid; h1 = findobj(gcf,'type','line'); set(h1,'linewidth',2);
19 subplot(224); step(P, P_delay, P_RHP_zero, P_RHP_pade); grid; h1 =
    findobj(gcf,'type','line'); set(h1,'linewidth',2);
20
21 legend

```

The step response of a RHP zero initially moves to the wrong direction, because the step input consists of both low and high frequency components (see Appendix ??). The RHP zero system will respond to high frequency input with much higher phase delay, and especially those near 180° will cause the system response to move to the wrong direction.

9.6. Rhor's Counter Example

Consider the plant transfer function

$$P_{\text{ideal}} = \frac{2}{s + 1} \quad (58)$$

We can design a closed loop system that has arbitrarily fast transients with a large enough gain K .

If the plant has some unmodelled dynamics, which could be true for example for a mechanical system has some elastic connector, the actual plant transfer function becomes

$$P_{\text{Rhor}} = \frac{2}{s + 1} \frac{229}{s^2 + 30s + 229} \quad (59)$$

Show that this system could become unstable when K is too large; and also show that if a constant gain controller is used, there is a maximum value for K that makes closed loop system stable.

10. Lead-Lag Compensator

We have learned root locus method and we will now see how it is useful in understanding the effect of lead-lag compensator³⁰ on the closed loop system roots. (TODO: adding an example rlocus plot.)

In chapter 7, we have learned that the fundamental guideline of transfer function based controller design is to make the open loop gain crossover occur at a slope of -20 dB/dec. We will soon find that this is a generally true:

- The lead compensator introduces the -20 dB/dec segment even when the plant is a type-2 system that only has slope of -40 dB/dec or lower.
- When the gain K has been determined to meet the steady state error constant requirement, the lag compensator can be used to modify the gain crossover frequency as an effort to move it to the -20 dB/dec segment of the magnitude-frequency curve. This requires the plant has a segment in Bode plot that already has a slope of -20 dB/dec. For example, system $P(s) = 1/s^2$ cannot be made stable by lag compensator.

10.1. Overview

Equivalent statement and conventions.

- Zero means open loop zero.
- Pole means open loop pole.
- Root means closed loop pole.
- Dominant root refers to dominant pole of closed loop system, because we are not concerned with the dominance of the open loop poles.
- A second order open loop transfer function with a marginally stable pole is the same as a first order system that is specified to track a ramp input with finite steady state error.

Fundamentally, there are five combinations in terms of design tools and design objectives.

- Bode plot can be used to estimate the phase margin introduced by the phase lead compensator, which moves the gain crossover frequency to the right of the uncompensated cut-off frequency. Lead compensator works even if there is no -20 dB/dec segment in the original Bode plot, because it introduces such slope near gain crossover.
- Bode plot aided design of phase lag compensator can also introduce phase margin by moving the gain crossover frequency to the left of the uncompensated cut-off frequency. This works only if there is already a -20 dB/dec segment in the original Bode plot.
- Root locus design of phase lead compensator basically puts a zero to the left of the dominant roots for adding damping (i.e., increasing damping ratio ζ).
- Root locus design of phase lag compensator focuses on steady state error reduction without changing the system type. This is achieved by adding a pair of real pole and zero near the origin, and their distance to the dominant roots must be sufficient.

³⁰In this course, compensator and controller are the same thing and there is no difference between the two. In most cases, the compensator is placed in the forward path, which is known as cascaded compensator as it is cascaded to the plant.

- Integrator is used to modify the system type, in order to meet the steady state error constant requirement. One integrator costs 90 degrees of phase margin, while lag compensator costs much less.

Relation to PD and PI controller.

- PD is an extreme case of phase lead compensator with the pole being placed at infinite.
- PI is an extreme case of phase lag compensator with the zero being placed at origin.

10.2. Phase Lead

Lead compensator can be used to add **damping** to dominant roots in root locus plot, and add **phase margin** in Bode plot. Also, the use of a lead compensator extends the **bandwidth** of a feedback system.

Root locus is a very satisfactory approach when the specifications are given in terms of percent overshoot and settling time, thus specifying the ζ and ω_n of the desired dominant roots in the s-plane [1].

When the design specifications include an error constant requirement, the Bode plot method is more suitable, because the root locus method often results in an iterative design procedure when the error constant is specified [1].

10.2.1. Lead Compensator Design Using Bode Plot

The peak phase margin can be calculated using the ratio between the pole and zero.

10.2.2. Lead Compensator Design Using Root Locus

Knowing the relation between phase margin and damping ratio of a standard second order system would be very helpful. Generally, it is suggested to place the zero below the dominant roots and to the left of the open loop poles.

The location of the added zero can be calculated by ϕ_1 and ϕ_2 . See Fig. ??, the location of the zero can be determined by

10.3. Phase Lag

Phase lag compensator introduces gain attenuation, which reduces the gain crossover frequency. Lag compensator also introduces unwanted phase lag between zero frequency and pole frequency, and therefore, we need to place it far away from the gain crossover region. It is generally recommended to place the compensator zero at 1/50 of the real part of the dominant poles.

10.3.1. Lag Compensator Design Using Bode Plot

Phase lag compensator can also be used to improve the phase margin, and it is achieved by advancing the gain crossover frequency.

10.3.2. Lag Compensator Design Using Root Locus

Note that the steady state error depends on the ratio of the zero location and pole location. By putting the zero and pole of the lag compensator far away from the

dominant roots, the lag compensator can be used to reduce the steady state error while keeping the performance of the dominant pole unchanged.

By adding a pair of pole and zero near imaginary axis and far away from the dominant poles, we can reduce steady state error and at the same time, ensure the impulse response is not changed. But, how about the phase margin, does it change?

10.3.3. Lead-Lag Compensator Design

- Plot the Bode plot of the open loop transfer function.
- Calculate for phase margin, and calculate for $\alpha = \frac{z_0}{p_0}$ using formula. If you forget about the formula, just try to use $\alpha = 10$. A lead compensator having α larger than 10 is not efficient in terms getting extra phase margin.
- Find the $-10 \log_{10}(\alpha)$ dB or -10 dB for $\alpha = 10$ frequency on the Bode plot,
- When the magnetite response of Bode plot crosses the 0 dB line with a slope of -40 dB/dec, a lag compensator can help to advance the gain crossover frequency, as an effort to increase the phase margin. This is because crossing 0 dB line with a slope of -40 dB/dec will “waste” phase margin a lot.
- The step response's steady state error with a compensator is calculated as follows.

$$\begin{aligned} L(s) &= C(s)P(s) = \frac{N(s)}{D(s)} \\ \Rightarrow E(s) &= R(s) - Y(s) = R(s) - \frac{L(s)}{1 + L(s)}R(s) \\ \Rightarrow e_{ss} &= \lim_{s \rightarrow 0} sE(s) = \frac{D(0)}{D(0) + KN(0)\frac{z}{p}} \end{aligned}$$

where $R(s) = 1/s$.

What's the gain? Where to put?

For lead compensator, the gain is alpha=10, where to put is the -10 dB of the Bode plot.

For lag compensator, the gain is ? put near the imaginary axis.

As seen from root locus plot, the frequency to place the lag compensator needs to avoid affecting the dominant pole

10.4. Two Different Kinds of Steady State Error

Four step design procedure when the magnitude gain at the required corner frequency is higher than -10 dB.

- (1) Set K to meet the speed error constant requirement $K_v = K$.
- (2) Design a lag compensator that provides a decrease in magnitude gain at required corner frequency ω_c^* : M_{lag}
- (3) Lead compensator is added to get extra phase margin.
- (4) Check corner frequency and phase margin of the open loop system after adding the compensator.

10.5. Two Different Kinds of Steady State Error

Consider a typical transfer function of a motor with neglectable viscous term.

The integral term contributed by the plant only makes the open loop transfer function for reference tracking a type 1 system.

Having an integral term in the plant does not make open loop transfer function for disturbance rejection a type 1 system. The steady state error caused by the disturbance can be derived as follows

a

In order to have a zero steady state error by disturbance, the controller must incorporate an integral term. This results in a type 2 system for reference tracking.

A type 2 system's phase-frequency response begins at -180° . If there is no zero added in the controller, its phase margin would be less or equal to zero.

This fact implies that the closed loop control system's stability and performance highly depends on the local characteristics of the Bode plot near the 0 dB gain crossing frequency.

10.6. PID Control in a Control Theory Perspective

Each term.

- The proportional term is the essence of the feedback control.
- The derivative term can be understood as adding an open loop zero that improves the damping of the closed loop system.
- The integral term puts the transfer function from step disturbance to error, $P(s)S(s)$, to zero with dc input.

Combinations.

- PD compensator introduces a zero (which introduces phase lead hence gain margin) but is not proper, Thus a high frequency pole should be added. The resulting compensator is the same as a lead compensator.
- PI compensator is able to remove steady state error under dc disturbance and at the same time, it introduces of a zero.
- PID compensator adds two zeros. To make it proper, two poles are added, one at origin and the other is placed at a higher frequency.

Arbitrary pole-placement can be done when the used controllers are of the same order as that of the plant.

Example of application of PID regulator. Design a compensator for $P(s) = 1/(s+1)$ that has zero steady state error under ramp input.

10.6.1. Practical Derivative Control

In practice, we need to add a low-pass filter derivative controller $C(s) = K_D s$ to have a proper controller that is physically realizable:

$$C(s) = \frac{1}{\tau s + 1} \times K_D s = \frac{K_D / \tau}{1 + \frac{1}{\tau s}} \quad (60)$$

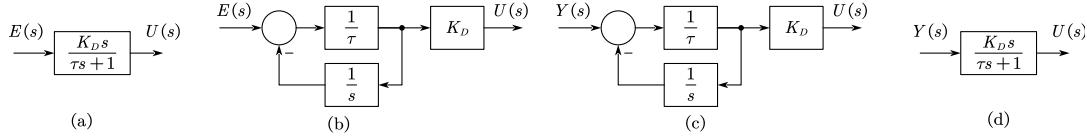


Figure 41. The proper derivative control implemented as a closed loop system. (a) proper derivative; (b) equivalent closed loop implementation; (c)(d) derivative-on-measurement.

where τ is the time constant of the first order low pass filter, and the second equal sign shows that this proper derivative control is equivalent to a closed loop system with an integrator in the feedback path, as shown in Fig. 41. The system in Fig. 41 can be implemented using fundamental blocks that might help algorithm efficiency.

It becomes apparent from Fig. 41 that a proper derivative or filtered derivative has an internal state, because an integrator introduces a state. In contrast, the ideal derivative control does not need a state and the calculation does not rely on information in the past.

10.6.2. Practical Integral Control

In practice, a switch needs to be added to the integral action. The integral action should be turned off when the actuator is saturated, when the actuator cannot respond to a too large command and its response winds up regardless of ever-increasing command. This turning off of integral action at the right moment is also known as **anti-windup**. Anti-windup can be implemented via dynamic clamping, back calculation and observer. We will now briefly provide two sample codes of the clamping method.

In Chapter 11 of TI InstaSPIN user's guide, there is one section about the dynamic anti-windup of the current controller. I do not recommend this implementation, but it is still useful to demonstrate the basic idea of anti-windup. Here is a code snippet that implements the PID with anti-windup mechanism using Tustin's method (ode2).

```

1 @njit(nogil=True)
2 def tustin_pid(reg):
3
4     # Error signal
5     error = reg.setpoint - reg.measurement
6
7     # Proportional
8     proportional = reg.Kp * error
9
10    # Integral
11    reg.integrator = reg.integrator + 0.5 * reg.Ki * reg.T * (error
12        + reg.prevError) # Tustin
13    # reg.integrator = reg.integrator + reg.Ki * reg.T * (error) # Euler
14
15    # Anti-wind-up via integrator clamping
16    if reg.integrator > reg.IntLimit:
17        reg.integrator = reg.IntLimit
18    elif reg.integrator < -reg.IntLimit:
19        reg.integrator = -reg.IntLimit
20
21    # Derivative (band-limited differentiator) # Note: derivative
22        # on measurement, therefore minus sign in front of equation!

```

```

1      */
2      reg.differentiator = -(2.0 * reg.Kd * (reg.measurement - reg.
3          prevMeasurement) \
4              + (2.0 * reg.tau - reg.T) * reg.
5                  differentiator) \
6              / (2.0 * reg.tau + reg.T)
7
8      # Compute output and apply limits
9      reg.Out = proportional + reg.integrator + reg.differentiator
10
11     if reg.Out > reg.OutLimit:
12         reg.Out = reg.OutLimit
13     elif reg.Out < -reg.OutLimit:
14         reg.Out = -reg.OutLimit
15
16     # Store error and measurement for later use
17     reg.prevError = error
18     reg.prevMeasurement = reg.measurement
19
20     # Implement dynamic clamping
21     reg.IntLimit = reg.OutLimit - proportional
22
23     # Return controller output
24     return reg.Out

```

Let's implement an ode1 version for incremental PI, in which the dynamic clamping is naturally implemented with only one block of output saturation. Note the sampling time has been absorbed into the PI gains.

```

1 @njit(nogil=True)
2 def incremental_pi(reg):
3     reg.Err = reg.setpoint - reg.measurement
4     reg.Out = reg.OutPrev + \
5         reg.Kp * (reg.Err - reg.ErrPrev) + \
6         reg.Ki * reg.Err
7     if reg.Out > reg.OutLimit:
8         reg.Out = reg.OutLimit
9     elif reg.Out < -reg.OutLimit:
10        reg.Out = -reg.OutLimit
11     reg.ErrPrev = reg.Err
12     reg.OutPrev = reg.Out

```

from the code, the following derivation can now be put forth

$$\begin{aligned}
u[k] - u[k-1] &= K_P(e[k] - e[k-1]) + K_I e[k] \\
\Rightarrow (1 - z^{-1}) u[k] &= K_P (1 - z^{-1}) e[k] + K_I e[k] \\
\Rightarrow u[k] &= K_P e[k] + \frac{K_I e[k]}{(1 - z^{-1})}
\end{aligned}$$

ODE1 means $s = \frac{1}{1-z^{-1}}$, which does not have a constant phase lag of 90° like ODE2 does.

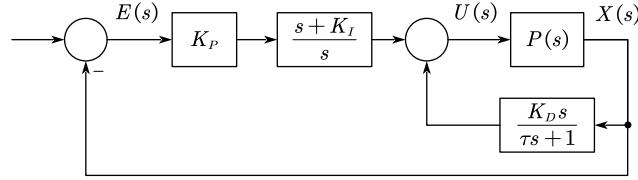


Figure 42. Block diagram of the PID controller that implements series PI and derivative-on-measurement.

10.6.3. Put'em Together

11. State Space Model

11.1. Conversion from Transfer Function

11.1.1. Controllable Canonical Form

11.1.2. Diagonal Canonical Form

11.1.3. The Duality

11.2. Controllability and Observability

A system is controllable, if every point in the state space is reachable. A point x^* is reachable, if there exists a solution v for equation $P_c v = x^*$. The span of the column of P_c defines the reachable subspace.

To determine which state is not controllable, we can find an invertible matrix P to make $PB = [0, 0, 0, 1]$. The controllability of the transformed system $\{PAP^{-1}, PB\}$ is equivalent to that of (A, B) . It is easier to only transform B into controllable canonical form.

11.3. Luenberger Observer

11.4. Separation Principle

11.5. Time Domain Solution of the State Space Model

The proof needs Leibniz integral rule.

11.6. Reference Input

The output needs to track a step input.

$$\begin{aligned}
 Y(s) &= C(sI - A - BK)^{-1} B \frac{N}{s} \\
 \lim_{s \rightarrow 0} sY(s) &= C(sI - A - BK)^{-1} BN = 1 \\
 \Rightarrow N &= \frac{1}{C(sI - A - BK)^{-1} B}
 \end{aligned} \tag{61}$$

which means a proper N value helps to eliminate steady state error of a step input.

Appendix A. Which to Use, MATLAB or Python?

[Important!] The code snippet in this document has preceding line numbers. It is possible to avoid selecting the line numbers.

- Holding ALT, and use box select in Adobe Acrobat Reader.
- Holding CTRL+ALT, and use box select in SumatraPDF.

MATLAB is friendly for beginners for having complete help documents and nice IDE. The use case to choose python over MATLAB is when I need to produce professional figure or make GUI. I am not suggesting there is no way to produce nice looking figure using MATLAB,³¹ but MATLAB is not a programming language and its OOP support is not as intuitive as others. There is an MATLAB computability mode for the python-control package.³² This course will use python package DearPyGUI for interactive learning if needed. DearPyGUI can take advantage of GPU and does not need a GUI editor like PyQt6 or PySide2.

Appendix B. Review Math Concepts: Two Kernels

A brief refresh on the concepts of kernels from the two major math courses is sufficient.

B.1. Kernel in Integral Transform

The integral transform is a math operation that changes variable of interest, and it has a general form as follows

$$F(\alpha) = \int_a^b f(t) K(\alpha, t) dt \quad (\text{B1})$$

where $K(\alpha, t)$ is known as the kernel of the integral transform.

There are three integral transforms will be used in this course: Laplace transform, Fourier transform and convolution:

$$F(s) = \int_{0^-}^{\infty} f(t) e^{-st} dt, \quad s \in \mathbb{C} \quad (\text{B2a})$$

$$F(j\omega) = \int_{0^-}^{\infty} f(t) e^{-j\omega t} dt, \quad \omega \in \mathbb{R} \quad (\text{B2b})$$

$$F(t) = \int_{0^-}^{\infty} f(\tau) \delta(t - \tau) d\tau, \quad \tau \in \mathbb{R} \quad (\text{B2c})$$

where $j = \sqrt{-1}$; $\delta(\cdot)$ is the impulse function or Dirac delta function that is nonzero only when $t \in [0^-, 0^+]$, the 0's superscript sign indicates one-sided limit—minus/plus sign indicates the limit approaches 0 from left/right side; signal $f(t)$ is assumed to be a causal signal, i.e., $f(t) = 0$ when $t < 0^-$, and this is why the integral begins at $t = 0^-$.

³¹See e.g., https://github.com/adinatan/plot_darkmode/tree/main

³²See ReadTheDocs <https://python-control.readthedocs.io/en/latest/matlab.html> or Source Forge https://python-control.sourceforge.net/manual/matlab_strings.html

B.2. Kernel in Linear Algebra

The rank–nullity theorem states that the number of columns of a matrix A is the sum of the rank of A and the nullity of A . This course is mainly concerned with the real matrix, so we can re-state rank–nullity theorem as follows

$$\text{rank}(A) + \text{nullity}(A) = n, \quad A \in \mathbb{R}^{m \times n} \quad (\text{B3})$$

where the nullity is the dimension of kernel. The kernel of a linear transform A is the space (i.e., a set of points) that is going to be mapped into a point (i.e., the origin):

$$\ker(A) = \{x \in \mathbb{R}^n \mid Ax = 0\} \quad (\text{B4})$$

and therefore, kernel is also known as the null space of matrix A . Looking at the definition of kernel (B4), we realize that

$$Av = 0 = 0v$$

which suggests matrix A has at least one eigenvalue that equals to 0.

In the sequel, we will assume A is a square n -by- n matrix, i.e., $m = n$.

B.2.1. Revisit Linear Map

Consider a \mathbb{R}^n space that is spanned by a series of n basis vectors:

A matrix is the mapping of the basis vectors.

The determinant is the length, area, and volume of the vectors.

Some linear map loses information during the mapping. Those information is mapped into a null space.

B.2.2. Eigenvalue and Eigenvector

The eigenvalue $\lambda \in \mathbb{C}$ of a matrix A satisfies:

$$Av = \lambda v \quad (\text{B5})$$

in which the vector $v \in \mathbb{R}^n$ is known as eigenvector. When $\lambda \in \mathbb{R}$, an eigenvector v only experiences a scale of λ after being transformed by matrix A . The eigenvectors can be used to span a space, known as eigenspace, and the following matrix diagonalization becomes possible:

$$\begin{aligned} [v_1, v_2, \dots, v_n] A &= [v_1, v_2, \dots, v_n] \begin{bmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \lambda_n \end{bmatrix} \\ \Rightarrow A &= [v_1, v_2, \dots, v_n] \begin{bmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \lambda_n \end{bmatrix} [v_1, v_2, \dots, v_n]^{-1} \end{aligned} \quad (\text{B6})$$

where v_i , $i = 1, 2, \dots, n$ is eigenvector.

There are two matrix properties related to eigenvalues³³

- The trace of a matrix $\text{tr}A$ equals to the sum of eigenvalues.

$$\text{tr}A = \sum_i^n \lambda_i \quad (\text{B7})$$

- The determinant of a matrix $\det A$ equals to the product of eigenvalues.

$$\det A = \prod_i^n \lambda_i \quad (\text{B8})$$

By finding the kernel or nullspace of the matrix $(A - \lambda I)$,³⁴ we can solve for the eigenvalues of matrix A :

$$(A - \lambda I) v = 0 = 0v \quad (\text{B9})$$

which suggests that 0 is an eigenvalue of matrix $(A - \lambda I)$. Using (B8), the above linear equations (B9) can be turned into a polynomial equation using the determinant of matrix A :

$$\det(A - \lambda I) = 0$$

which gives the eigenvalues of matrix A .

The trace of a matrix A is the divergence of this vector field created by the matrix A .

B.2.3. Eigenvalue and Eigenfunction

There is an analogy between eigenvector and exponential mode, if we replace the linear operator A defined in a space of finite dimension with derivative $\frac{d}{dt}$ defined in a space of infinite dimension. The “eigenvector” of operator $\frac{d}{dt}$ is the exponential function.³⁵ This can be shown by solving the following eigenvalue equation:

$$\begin{aligned} \frac{d}{dt} f(t) &= \lambda f(t) \\ \Rightarrow \frac{df(t)}{f(t)} &= \lambda dt \\ \Rightarrow \ln |f(t) - f(0)| &= \lambda t \\ \Rightarrow f(t) &= f(0)e^{\lambda t} \end{aligned} \quad (\text{B10})$$

in which $e^{\lambda t}$ is formally known as the eigenfunction. The exponential modes can be used to describe the solution of O.D.E.

³³https://en.wikipedia.org/wiki/Eigenvalues_and_eigenvectors#Additional_properties_of_eigenvalues

³⁴https://en.wikipedia.org/wiki/Eigenvalues_and_eigenvectors#Eigenspaces,_geometric_multiplicity,_and_the_eigenbasis_for_matrices

³⁵https://en.wikipedia.org/wiki/Eigenvalues_and_eigenvectors#Eigenvalues_and_eigenfunctions_of_differential_operators

B.2.4. Matrix Exponential

For a square matrix A , the matrix exponential is defined in terms of the following “Taylor series”-alike matrix power series:³⁶

$$\begin{aligned} e^{At} &= \sum_{k=0}^{\infty} \frac{t^k A^k}{k!} = I + At + \frac{AA}{2!}t^2 + \frac{AAA}{3!}t^3 + \dots \\ &= [v_1, v_2, \dots, v_n] \exp \left(\begin{bmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \lambda_n \end{bmatrix} t \right) [v_1, v_2, \dots, v_n]^{-1} \end{aligned} \quad (\text{B11})$$

As a result, the solution to the following state space model

$$\frac{d}{dt}x = Ax$$

becomes

$$\begin{aligned} x(t) &= x(0)e^{At} \\ &= x(0)[v_1, v_2, \dots, v_n] \exp \left(\begin{bmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \lambda_n \end{bmatrix} t \right) [v_1, v_2, \dots, v_n]^{-1} \end{aligned}$$

where note x is the vector of state variables, and v_i is eigenvector, and they are related by the homogeneous matrix A .

In addition, the following equality holds for determinant of matrix exponential:

$$e^{\text{tr}A} = \det(e^A) \neq 0$$

Appendix C. Five Ways Solving Ordinary Differential Equations

See “Five Levels for Differential Equations in Physics” by Physics with Elliot <https://www.physicswithelliot.com/odes-help-room-notes>

Appendix D. Zeros and Zero dynamics

Appendix E. Passivity and Stability Margin

The stability margin of a control system is reduced when there are unmodelled dynamics in the plant. See video Brian’s video “Passivity-Based Control to Guarantee Stability”, where the passivity based tuning guideline has been used to improve the relative stability of a LQG controller.

The positive real (PR) condition needed in the KYP lemma³⁷ is related to passivity. There is a dedicated book written by Remeo Ortega for passivity based control.

³⁶<https://zhuanlan.zhihu.com/p/57051153>

³⁷https://en.wikipedia.org/wiki/Kalman%20%20%20%20Yakubovich%20%20%20Popov_lemma

Appendix F. Fourier Analysis and Time-Frequency Domain Analysis

F.1. Frequency Response and Fourier Analysis

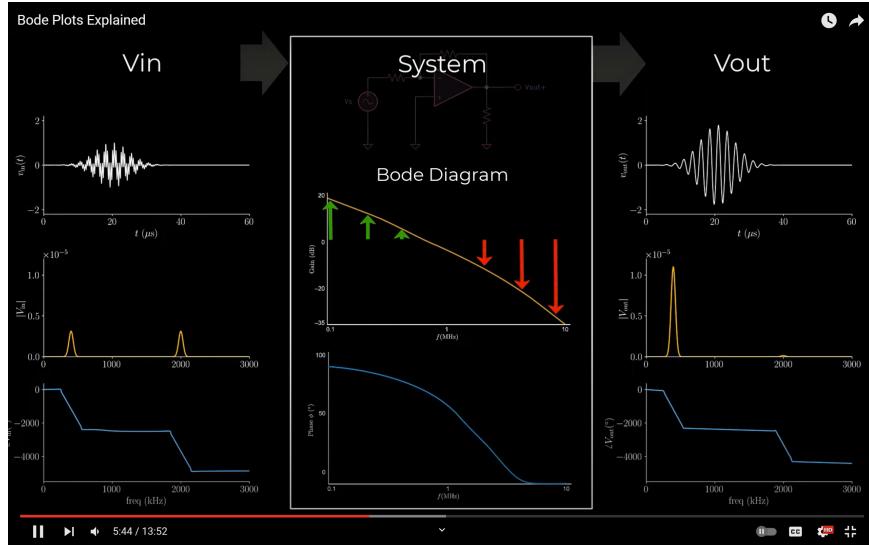


Figure F1. The difference between Bode plot and Fourier analysis. Screenshot from Youtube video: PF4fSRwPk5I.

The Bode plot is about relation between input and output signal. The Fourier analysis is concerned with signal of certain periods. See Fig. F1.

F.2. Time Frequency Domain Analysis

The DFT works fine for a non-periodic signal such as Heaviside step function, in which only the frequencies of the bins are reliable, while the magnitudes of the bins are dependent on the signal duration.

A better tool for analyzing non-periodic rich frequencies signal is the time-frequency analysis. Here is a sample code.

```

1  %% Time-Frequency Analysis
2  clear
3  clc
4
5  Fs = 1e3 % Sampling frequency
6
7  time = 0:1/Fs:10-1/Fs; % make length an even number of points for
   fft
8  a = 0*time;
9  b = sin(2*pi*time).*cos(2*pi*100.*time) + sin(2*pi*10*time);
10 b = sign(time);
11 c = 0*time;
12 signal = [a b c];
13 subplot(131)
14 pspectrum(signal, Fs, 'spectrogram', 'Reassign', true, ...
15   'FrequencyLimits', [0 150], 'TimeResolution', 1.0)
16
17
18 T = 1/Fs; % Sampling period

```

```

19 L = length(b); % Length of signal
20
21 % t = (0:L-1)*T; % Time vector
22 % S = 0.7*sin(2*pi*50*t) + sin(2*pi*120*t);
23 % X = S + 2*randn(size(t));
24 %
25
26 subplot(132)
27 plot(signal)
28 title("Signal Corrupted with Zero-Mean Random Noise")
29 xlabel("t (milliseconds)")
30 ylabel("X(t)")
31
32 Y = fft(b);
33 Y = fft(signal);
34
35 P2 = abs(Y/L);
36 P1 = P2(1:L/2+1);
37 P1(2:end-1) = 2*P1(2:end-1);
38
39 f = Fs*(0:(L/2))/L;
40 subplot(133)
41 plot(f, P1, '.')
42 title("Single-Sided Amplitude Spectrum of X(t)")
43 xlabel("f (Hz)")
44 ylabel("|P1(f)|")

```

Appendix G. Cascaded Loop Control

Cascaded loop control is useful to separate the problem.

Appendix H. Pending Proofs

As a fundamental course in the field of control, there are lacking proofs for a few established results.

- Routh-Hurwitz Criterion?
- Why frequency response can be obtained by substituting $s = j\omega$?
- Why does contour mapping work?
- Argument principle.

Appendix I. Textbook Errata

The following errata can be made to improve [1].

- The Laplace transform table in the appendix is not 100% correct.
- The pole should be $-p$ and zero should be $-z$ —the minus sign is occasionally missing in the text.
- The definition of nonminimum phase system is not complete.

References

- [1] R. C. Dorf and R. H. Bishop, *Modern Control Systems (14th Ed.)*. Pearson Education Limited, 2020.
- [2] B. Douglas, *The Fundamentals of Control Theory (Rev 1.6)*. EngineeringMedia, 2019.
- [3] G. C. Goodwin, S. F. Graebe, and M. E. Salgado, *Control System Design*. Pearson, 2000.
- [4] H. K. Khalil, *Control Systems: An Introduction*, 2023.
- [5] K. Ogata, *Modern Control Engineering 5th ed.* Prentice Hall, 2010.