

SECTION 19

CONTROL SYSTEMS

Control is used to modify the behavior of a system so it behaves in a specific desirable way over time. For example, we may want the speed of a car on the highway to remain as close as possible to 60 miles per hour in spite of possible hills or adverse wind; or we may want an aircraft to follow a desired altitude, heading, and velocity profile independent of wind gusts; or we may want the temperature and pressure in a reactor vessel in a chemical process plant to be maintained at desired levels. All these are being accomplished today by control methods and the above are examples of what automatic control systems are designed to do, without human intervention. Control is used whenever quantities such as speed, altitude, temperature, or voltage must be made to behave in some desirable way over time.

This section provides an introduction to control system design methods. P.A., Z.G.

In This Section:

CHAPTER 19.1 CONTROL SYSTEM DESIGN	19.3
INTRODUCTION	19.3
Proportional-Integral-Derivative Control	19.3
The Role of Control Theory	19.4
MATHEMATICAL DESCRIPTIONS	19.4
Linear Differential Equations	19.4
State Variable Descriptions	19.5
Transfer Functions	19.7
Frequency Response	19.9
ANALYSIS OF DYNAMICAL BEHAVIOR	19.10
System Response, Modes and Stability	19.10
Response of First and Second Order Systems	19.11
Transient Response Performance Specifications for a Second Order Underdamped System	19.13
Effect of Additional Poles and Zeros	19.14
CLASSICAL CONTROL DESIGN METHODS	19.14
Design Specifications and Constraints	19.14
Control Design Strategy Overview	19.15
Evaluation of Control System	19.19
Digital Implementation	19.20
ALTERNATIVE DESIGN METHODS	19.21
Nonlinear PID	19.21
State Feedback and Observer Based-Design	19.22

ADVANCED ANALYSIS AND DESIGN TECHNIQUES	19.26
APPENDIX: OPEN AND CLOSED LOOP STABILIZATION	19.27
REFERENCES	19.29



On the CD-ROM:

“A Brief Review of the Laplace Transform,” by the authors of this section, examines its usefulness in control Functions.

CHAPTER 19.1

CONTROL SYSTEM DESIGN

Panus Antsaklis, Zhiqian Gao

INTRODUCTION

To gain some insight into how an automatic control system operates we shall briefly examine the speed control mechanism in a car.

It is perhaps instructive to consider first how a typical driver may control the car speed over uneven terrain. The driver, by carefully observing the speedometer, and appropriately increasing or decreasing the fuel flow to the engine, using the gas pedal, can maintain the speed quite accurately. Higher accuracy can perhaps be achieved by looking ahead to anticipate road inclines. An automatic speed control system, also called *cruise control*, works by using the difference, or error, between the actual and desired speeds and knowledge of the car's response to fuel increases and decreases to calculate via some algorithm an appropriate gas pedal position, so to drive the speed error to zero. This decision process is called a *control law* and it is implemented in the *controller*. The system configuration is shown in Fig. 19.1.1. The car dynamics of interest are captured in the *plant*. Information about the actual speed is fed back to the controller by *sensors*, and the control decisions are implemented via a device, the *actuator*, that changes the position of the gas pedal. The knowledge of the car's response to fuel increases and decreases is most often captured in a mathematical model.

Certainly in an automobile today there are many more automatic control systems such as the antilock brake system (ABS), emission control, and tracking control. The use of feedback control preceded control theory, outlined in the following sections, by over 2000 years. The first feedback device on record is the famous Water Clock of Ktesibios in Alexandria, Egypt, from the third century BC.

Proportional-Integral-Derivative Control

The proportional-integral-derivative (PID) controller, defined by

$$u = K_p e + K_I \int e + K_D \dot{e} \quad (1)$$

is a particularly useful control approach that was invented over 80 years ago. Here K_p , K_I , and K_D are controller parameters to be selected, often by trial and error or by the use of a lookup table in industry practice. The goal, as in the cruise control example, is to drive the error to zero in a desirable manner. All three terms Eq. (1) have explicit physical meanings in that e is the current error, $\int e$ is the accumulated error, and \dot{e} represents the trend. This, together with the basic understanding of the causal relationship between the control signal (u) and the output (y), forms the basis for engineers to “tune,” or adjust, the controller parameters to meet the design specifications. This intuitive design, as it turns out, is sufficient for many control applications.

To this day, PID control is still the predominant method in industry and is found in over 95 percent of industrial applications. Its success can be attributed to the simplicity, efficiency, and effectiveness of this method.

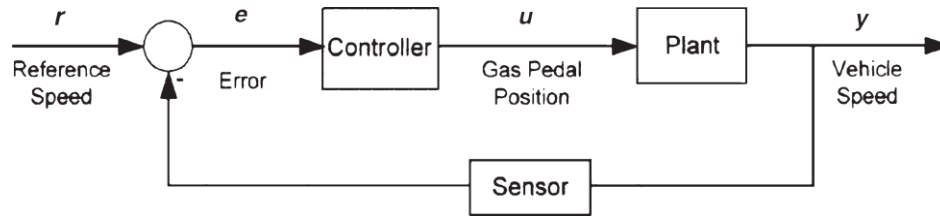


FIGURE 19.1.1 Feedback control configuration with cruise control as an example.

The Role of Control Theory

To design a controller that makes a system behave in a desirable manner, we need a way to predict the behavior of the quantities of interest over time, specifically how they change in response to different inputs. Mathematical models are most often used to predict future behavior, and control system design methodologies are based on such models. Understanding control theory requires engineers to be well versed in basic mathematical concepts and skills, such as solving differential equations and using Laplace transform. The role of control theory is to help us gain insight on how and why feedback control systems work and how to systematically deal with various design and analysis issues. Specifically, the following issues are of both practical importance and theoretical interest:

1. Stability and stability margins of closed-loop systems.
2. How fast and smooth the error between the output and the set point is driven to zero.
3. How well the control system handles unexpected external disturbances, sensor noises, and internal dynamic changes.

In the following, modeling and analysis are first introduced, followed by an overview of the classical design methods for single-input single-output plants, design evaluation methods, and implementation issues. Alternative design methods are then briefly presented. Finally, For the sake of simplicity and brevity, the discussion is restricted to linear, time invariant systems. Results maybe found in the literature for the cases of linear, time-varying systems, and also for nonlinear systems, systems with delays, systems described by partial differential equations, and so on; these results, however, tend to be more restricted and case dependent.

MATHEMATICAL DESCRIPTIONS

Mathematical models of physical processes are the foundations of control theory. The existing analysis and synthesis tools are all based on certain types of mathematical descriptions of the systems to be controlled, also called plants. Most require that the plants are linear, causal, and time invariant. Three different mathematical models for such plants, namely, linear ordinary differential equation, state variable or state space description, and transfer function are introduced below.

Linear Differential Equations

In control system design the most common mathematical models of the behavior of interest are, in the time domain, linear ordinary differential equations with constant coefficients, and in the frequency or transform domain, transfer functions obtained from time domain descriptions via Laplace transforms.

Mathematical models of dynamic processes are often derived using physical laws such as Newton's and Kirchhoff's. As an example consider first a simple mechanical system, a spring/mass/damper. It consists of a weight m on a spring with spring constant k , its motion damped by friction with coefficient f (Fig. 19.1.2).

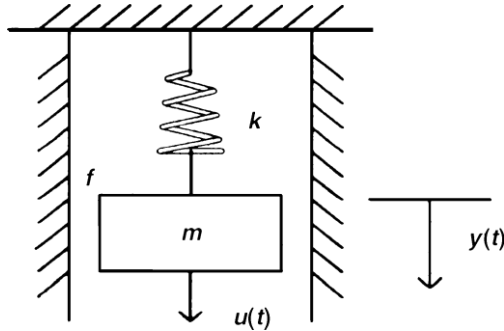


FIGURE 19.1.2 Spring, mass, and damper system.

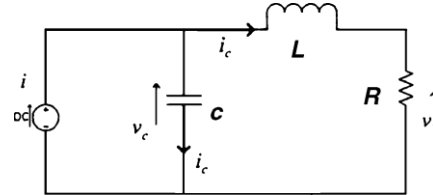


FIGURE 19.1.3 RLC circuit.

If $y(t)$ is the displacement from the resting position and $u(t)$ is the force applied, it can be shown using Newton's law that the motion is described by the following linear, ordinary differential equation with constant coefficients:

$$\ddot{y}(t) + \frac{f}{m} \dot{y}(t) + \frac{k}{m} y(t) = \frac{1}{m} u(t)$$

where $\dot{y}(t) \triangleq dy(t)/dt$ with initial conditions

$$y(t) \Big|_{t=0} = y(0) = y_0 \quad \text{and} \quad \frac{dy(t)}{dt} \Big|_{t=0} = \dot{y}(0) = \dot{y}_0$$

Note that in the next subsection the trajectory $y(t)$ is determined, in terms of the system parameters, the initial conditions, and the applied input force $u(t)$, using a methodology based on Laplace transform. The Laplace transform is briefly reviewed in Appendix A.

For a second example consider an electric RLC circuit with $i(t)$ the input current of a current source, and $v(t)$ the output voltage across a load resistance R . (Fig. 19.1.3)

Using Kirchhoff's laws one may derive:

$$\ddot{v}(t) + \frac{R}{L} \dot{v}(t) + \frac{1}{LC} v(t) = \frac{R}{LC} i(t)$$

which describes the dependence of the output voltage $v(t)$ to the input current $i(t)$. Given $i(t)$ for $t \geq 0$, the initial values $v(0)$ and $\dot{v}(0)$ must also be given to uniquely define $v(t)$ for $t \geq 0$.

It is important to note the similarity between the two differential equations that describe the behavior of a mechanical and an electrical system, respectively. Although the interpretation of the variables is completely different, their relations described by the same linear, second-order differential equation with constant coefficients. This fact is well understood and leads to the study of mechanical, thermal, fluid systems via convenient electric circuits.

State Variable Descriptions

Instead of working with many different types of higher-order differential equations that describe the behavior of the system, it is possible to work with an equivalent set of standardized first-order vector differential equations that can be derived in a systematic way. To illustrate, consider the spring/mass/damper example. Let $x_1(t) = y(t)$, $x_2(t) = \dot{y}(t)$ be new variables, called *state variables*. Then the system is equivalently described by the equations

$$\begin{aligned} \dot{x}_1(t) &= x_2(t) \quad \text{and} \quad \dot{x}_2(t) = -\frac{f}{m} x_2(t) - \frac{k}{m} x_1(t) + \frac{1}{m} u(t) \end{aligned}$$

with initial conditions $x_1(0) = y_0$ and $x_2(0) = y_1$. Since $y(t)$ is of interest, the output equation $y(t) = x_1(t)$ is also added. These can be written as

$$\begin{bmatrix} \dot{x}_1(t) \\ \dot{x}_2(t) \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -k/m & -f/m \end{bmatrix} \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} + \begin{bmatrix} 0 \\ 1/m \end{bmatrix} u(t)$$

$$y(t) = \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix}$$

which are of the general form

$$\dot{x}(t) = Ax(t) + Bu(t), \quad y(t) = Cx(t) + Du(t)$$

Here $x(t)$ is a 2×1 vector (a column vector) with elements the two state variables $x_1(t)$ and $x_2(t)$. It is called the *state vector*. The variable $u(t)$ is the *input* and $y(t)$ is the *output* of the system. The first equation is a vector differential equation called the *state equation*. The second equation is an algebraic equation called the *output equation*. In the above example $D = 0$; D is called the direct link, as it directly connects the input to the output, as opposed to connecting through $x(t)$ and the dynamics of the system. The above description is the *state variable or state space description* of the system. The advantage is that, system descriptions can be written in a standard form (the state space form) for which many mathematical results exist. We shall present a number of them in this section.

A state variable description of a system can sometimes be derived directly, and not through a higher-order differential equation. To illustrate, consider the circuit example presented above: using Kirchhoff's current law

$$i_c = C \frac{dv_c}{dt} = i - i_L$$

and from the voltage law

$$L \frac{di_L}{dt} = -Ri_L + v_c$$

If the state variables are selected to be $x_1 = v_c$, $x_2 = i_L$, then the equations may be written as

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & -1/C \\ 1/L & -R/L \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 1/C \\ 0 \end{bmatrix} v$$

$$v = \begin{bmatrix} 0 & R \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

where $v = Ri_L = Rx_2$ is the output of interest. Note that the choice of state variables is not unique. In fact, if we start from the second-order differential equation and set $\bar{x}_1 = v$ and $\bar{x}_2 = v'$, we derive an equivalent state variable description, namely,

$$\begin{bmatrix} \dot{\bar{x}}_1 \\ \dot{\bar{x}}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -1/LC & -R/L \end{bmatrix} \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \end{bmatrix} + \begin{bmatrix} 0 \\ R/LC \end{bmatrix} u$$

$$v = \begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \end{bmatrix}$$

Equivalent state variable descriptions are obtained by a change in the basis (coordinate system) of the vector state space. Any two equivalent representations

$$x = Ax + Bu, \quad y = Cx + Du \quad \text{and} \quad \bar{x} = A\bar{x} + \bar{B}u, \quad y = \bar{C}\bar{x} + \bar{D}u$$

are related by $\bar{A} = PAP^{-1}$, $\bar{B} = PB$, $\bar{C} = CP^{-1}$, $\bar{D} = D$, and $\bar{x} = Px$ where P is a square and nonsingular matrix. Note that state variables can represent physical quantities that may be measured, for instance, $x_1 = v_c$ voltage, $x_2 = i_L$ current in the above example; or they can be mathematical quantities, which may not have direct physical interpretation.

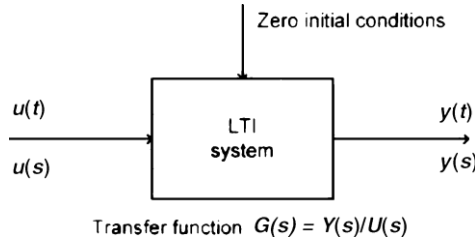


FIGURE 19.1.4 The transfer function model.

Linearization. The linear models studied here are very useful not only because they describe linear dynamical processes, but also because they can be approximations of nonlinear dynamical processes in the neighborhood of an operating point. The idea in linear approximations of nonlinear dynamics is analogous to using Taylor series approximations of functions to extract a linear approximation. A simple example is that of a simple pendulum $\dot{x}_1 = x_2$, $\dot{x}_2 = -k \sin x_1$, where for small excursions from the equilibrium at zero, $\sin x_1$ is approximately equal to x_1 and the equations become linear, namely, $\dot{x}_1 = x_2$, $\dot{x}_2 = -kx_1$.

Transfer Functions

The *transfer function* of a linear, time-invariant system is the ratio of the Laplace transform of the output $Y(s)$ to the Laplace transform of the corresponding input $U(s)$ with all initial conditions assumed to be zero (Fig. 19.1.4).

From Differential Equations to Transfer Functions. Let the equation

$$\frac{d^2 y(t)}{dt^2} + a_1 \frac{dy(t)}{dt} + a_0 y(t) = b_0 u(t)$$

with some initial conditions

$$y(t)\Big|_{t=0} = y(0) \quad \text{and} \quad \frac{dy(t)}{dt}\Big|_{t=0} = \frac{dy(0)}{dt} = y'(0)$$

describe a process of interest, for example, a spring/mass/damper system; see previous subsection. Taking Laplace transform of both sides we obtain

$$[s^2 y(s) - sy(0) - y'(0)] + a_1 [sY(s) - y(0)] + a_0 Y(s) = b_0 U(s)$$

where $Y(s) = L\{y(t)\}$ and $U(s) = L\{u(t)\}$. Combining terms and solving with respect to $Y(s)$ we obtain:

$$Y(s) = \frac{b_0}{s^2 + a_1 s + a_0} U(s) + \frac{(s + a_1)y(0) + y'(0)}{s^2 + a_1 s + a_0}$$

Assuming the initial conditions are zero,

$$Y(s)/U(s) = G(s) = \frac{b_0}{s^2 + a_1 s + a_0}$$

where $G(s)$ is the transfer function of the system defined above.

We are concerned with transfer functions $G(s)$ that are rational functions, that is, ratios of polynomials in s , $G(s) = n(s)/d(s)$. We are interested in proper $G(s)$ where $\lim_{s \rightarrow \infty} G(s) < \infty$. Proper $G(s)$ have degree $n(s) \leq$ degree $d(s)$.

In most cases $\text{degree } n(s) < \text{degree } d(s)$, which case $G(s)$ is called *strictly proper*. Consider the transfer function

$$G(s) = \frac{b_m s^m + b_{m-1} s^{m-1} + \dots + b_1 s + b_0}{s^n + a_{n-1} s^{n-1} + \dots + a_1 s + a_0} \quad \text{with } m \leq n$$

Note that the system described by this $G(s)$ ($Y(s) = G(s)U(s)$) is described in the time domain by the following differential equation:

$$y^{(n)}(t) + a_{n-1}y^{(n-1)}(t) + \dots + a_1 y'(t) + a_0 y(t) = b_m u^{(m)}(t) + \dots + b_1 u'(t) + b_0 u(t)$$

where $y^{(n)}(t)$ denotes the n th derivative of $y(t)$ with respect to time t . Taking Laplace transform of both sides of this differential equation, assuming that all initial conditions are zero, one obtains the above transfer function $G(s)$.

From State Space Descriptions to Transfer Functions. Consider $\dot{x}(t) = Ax(t) + Bu(t)$, $y(t) = Cx(t) + Du(t)$ with $x(0)$ initial conditions; $x(t)$ is in general an n -tuple, that is a (column) vector with n elements. Taking Laplace transform of both sides of the state equation:

$$sX(s) - x(0) = AX(s) + BU(s) \text{ or } (sI_n - A)X(s) = BU(s) + x(0)$$

where I_n is the $n \times n$ identity matrix; it has 1 on all diagonal elements and 0 everywhere else, e.g.,

$$I_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

Then

$$X(s) = (sI_n - A)^{-1}BU(s) + (sI_n - A)^{-1}x(0)$$

Taking now Laplace transform on both sides of the output equation we obtain $Y(s) = CX(s) + DU(s)$. Substituting we obtain,

$$Y(s) = [C(sI_n - A)^{-1}B + D]U(s) + C(sI_n - A)^{-1}x(0)$$

The response $y(t)$ is the inverse Laplace of $Y(s)$. Note that the second term on the right-hand side of the expression depends on $x(0)$ and it is zero when the initial conditions are zero, i.e., when $x(0) = 0$. The first term describes the dependence of Y on U and it is not difficult to see that the transfer function $G(s)$ of the systems is

$$G(s) = C(sI_n - A)^{-1}B + D$$

Example Consider the spring/mass/damper example discussed previously with state variable description $\dot{x} = Ax + Bu$, $y = Cx$. If $m = 1$, $f = 3$, $k = 2$, then

$$A = \begin{bmatrix} 0 & 1 \\ -2 & -3 \end{bmatrix}, \quad B = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad C = [1 \ 0]$$

and its transfer function $G(s)$ ($Y(s) = G(s)U(s)$) is

$$\begin{aligned} G(s) &= C(sI_2 - A)^{-1}B = [1 \ 0] \begin{bmatrix} s & 0 \\ 0 & s \end{bmatrix}^{-1} \begin{bmatrix} 0 & 1 \\ -2 & -3 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \\ &= [1 \ 0] \begin{bmatrix} s & -1 \\ 2 & s+3 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = [1 \ 0] \begin{bmatrix} 1 \\ s^2 + 3s + 2 \end{bmatrix} \begin{bmatrix} s+3 & 1 \\ s & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \\ &= \frac{1}{s^2 + 3s + 2} \end{aligned}$$

as before.

Using the state space description and properties of Laplace transform an explicit expression for $y(t)$ in terms of $u(t)$ and $x(0)$ may be derived. To illustrate, consider the *scalar case* $\dot{z} = az + bu$ with $z(0)$ initial condition. Using Laplace transform:

$$Z(s) = \frac{1}{s-a} z(0) + \frac{b}{s-a} U(s)$$

from which

$$z(t) = L^{-1}\{Z(s)\} = e^{at} z(0) + \int_0^t e^{a(t-\tau)} bu(\tau) d\tau$$

Note that the second term is a convolution integral. Similarly in the *vector case*, given

$$\dot{x}(t) = Ax(t) + Bu(t), \quad y(t) = Cx(t) + B(t)u(t)$$

it can be shown that

$$x(t) = e^{At} x(0) + \int_0^t e^{A(t-\tau)} Bu(\tau) d\tau$$

and

$$y(t) = Ce^{At} x(0) + \int_0^t Ce^{A(t-\tau)} Bu(\tau) d\tau + Du(t)$$

Notice that $e^{At} = L^{-1}\{(sI - A)^{-1}\}$. The *matrix exponential* e^{At} is defined by the (convergent) series

$$e^{At} = I + e^{At} + \frac{A^2 t^2}{2!} + \frac{A^k t^k}{k!} + \dots = I + \sum_{k=1}^{\infty} \frac{t^k}{k!} A$$

Poles and Zeros. The n roots of the denominator polynomial $d(s)$ of $G(s)$ are the *poles* of $G(s)$. The m roots of the numerator polynomial $n(s)$ of $G(s)$ are (finite) *zeros* of $G(s)$.

Example (Fig. 19.1.5)

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

$G(s)$ has one (finite) zero at -2 and two complex conjugate poles at $-1 \pm j$

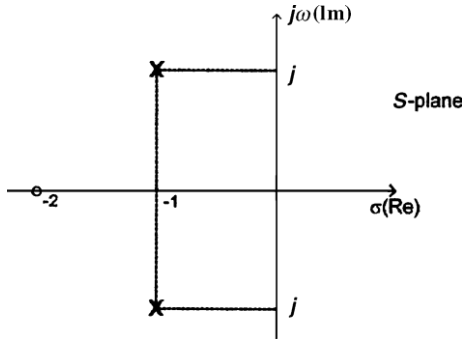
In general, a transfer function with m zeros and n poles can be written as

$$G(s) = k \frac{(s-z_1) \dots (s-z_m)}{(s-p_1) \dots (s-p_n)}$$

where k is the *gain*.

Frequency Response

The *frequency response* of a system is given by its transfer function $G(s)$ evaluated at $s = j\omega$, that is, $G(j\omega)$. The frequency response is a very useful means characterizing a system, since typically it can be determined experimentally, and since control system specifications are frequently expressed in terms of the frequency response. When the poles of $G(s)$ have negative real parts, the system turns out to be bounded-input/bounded-output

FIGURE 19.1.5 Complex conjugate poles of $G(s)$.

(BIBO) stable. Under these conditions the frequency response $G(j\omega)$ has a clear physical meaning, and this fact can be used to determine $G(j\omega)$ experimentally. In particular, it can be shown that if the input $u(t) = k \sin(\omega_o t)$ is applied to a system with a stable transfer function $G(s)$ ($Y(s) = G(s)U(s)$), then the output $y(t)$ at steady state (after all transients have died out) is given by

$$y_{ss}(t) = k |G(\omega_o)| \sin[\omega_o t + \theta(\omega_o)]$$

where $|G(\omega_o)|$ denotes the magnitude of $G(j\omega_o)$ and $\theta(\omega_o) = \arg G(j\omega_o)$ is the argument or phase of the complex quantity $G(j\omega_o)$. Applying sinusoidal inputs with different frequencies ω_o and measuring the magnitude and phase of the output at steady state, it is possible to determine the full frequency response of the system $G(j\omega_o) = |G(\omega_o)| e^{j\theta(\omega_o)}$.

ANALYSIS OF DYNAMICAL BEHAVIOR

System Response, Modes and Stability

It was shown above how the response of a system to an input and under some given initial conditions can be calculated from its differential equation description using Laplace transforms. Specifically, $y(t) = L^{-1}\{Y(s)\}$ where

$$Y(s) = \frac{n(s)}{d(s)}U(s) + \frac{m(s)}{d(s)}$$

with $n(s)/d(s) = G(s)$, the system transfer function; the numerator $m(s)$ of the second term depends on the initial conditions and it is zero when all initial conditions are zero, i.e., when the system is initially at rest.

In view now of the partial fraction expansion rules, see Appendix A, $Y(s)$ can be written as follows:

$$Y(s) = \frac{c_1}{s - p_1} + L + \frac{c_{i1}}{s - p_i} + \frac{c_{i2}}{(s - p_i)^2} + L + \frac{b_1 s + b_0}{s^2 + a_1 s + a_0} + L + I(s)$$

This expression shows real poles of $G(s)$, namely, p_1, p_2, \dots , and it allows for multiple poles p_i ; it also shows complex conjugate poles $a \pm jb$ written as second-order terms. $I(s)$ denotes the terms due to the input $U(s)$; they are fractions with poles the poles of $U(s)$. Note that if $G(s)$ and $U(s)$ have common poles they are combined to form multiple-pole terms.

Taking now the inverse Laplace transform of $Y(s)$:

$$y(t) = L^{-1}\{Y(s)\} = c_1 e^{p_1 t} + L + c_{i1} e^{p_i t} + (\cdot) t e^{p_i t} + L + e^{at}[(\cdot) \sin bt + (\cdot) \cos bt] + L + i(t)$$

where $i(t)$ depends on the input. Note that the terms of the form $c t^k e^{p t}$ are the *modes of the system*. The system behavior is the aggregate of the behaviors of the modes. Each mode depends primarily on the location of the pole p_i ; the location of the zeros affects the size of its coefficient c .

If the input $u(t)$ is a *bounded signal*, i.e., $|u(t)| < \infty$ for all t , then all the poles of $I(s)$ have real parts that are negative or zero, and this implies that $I(t)$ is also bounded for all t . In that case, the response $y(t)$ of the system will also be bounded for any bounded $u(t)$ if and only if all the poles of $G(s)$ have strictly negative real parts. Note that poles of $G(s)$ with real parts equal to zero are not allowed, since if $U(s)$ also has poles at the same locations, $y(t)$ will be unbounded. Take, for example, $G(s) = 1/s$ and consider the bounded step input $U(s) = 1/s$; the response $y(t) = t$, which is not bounded.

Having all the poles of $G(s)$ located in the open left half of the s -plane is very desirable and it corresponds to the system being stable. In fact, a system is *bounded-input, bounded-output (BIBO) stable* if and only if all poles of its transfer function have negative real parts. If at least one of the poles has positive

real part, then the system is *unstable*. If a pole has zero real part, then the term *marginally stable* is sometimes used.

Note that in a BIBO stable system if there is no forcing input, but only initial conditions are allowed to excite the system, then $y(t)$ will go to zero as t goes to infinity. This is a very desirable property for a system to have, because nonzero initial conditions always exist in most real systems. For example, disturbances such as interference may add charge to a capacitor in an electric circuit, or a sudden brief gust of wind may change the heading of an aircraft. In a stable system the effect of the disturbances will diminish and the system will return to its previous desirable operating condition. For these reasons a control system should first and foremost be guaranteed to be stable, that is, it should always have poles with negative real parts. There are many design methods to stabilize a system or if it is initially stable to preserve its stability, and several are discussed later in this section.

Response of First and Second Order Systems

Consider a system described by a first-order differential equation, namely, $\dot{y}(t) + a_0 y(t) = a_0 u(t)$ and let $y(0) = 0$. In view of the previous subsection, the transfer function of the system is

$$G(s) = \frac{a_0}{s + a_0}$$

and the response to a *unit step input* $q(t)$ ($q(t) = 1$ for $t \geq 0$, $q(t) = 0$ for $t < 0$) may be found as follows:

$$\begin{aligned} y(t) &= L^{-1}\{Y(s)\} = L^{-1}\{G(s)U(s)\} = L^{-1}\left\{\frac{a_0}{s+a_0} - \frac{1}{s}\right\} \\ &= L^{-1}\left\{\frac{1}{s} + \frac{-1}{s+a_0}\right\} = [1 - e^{-a_0 t}]q(t) \end{aligned}$$

Note that the pole of the system is $p = -a_0$ (in Fig. 19.1.6 we have assumed that $a_0 > 0$). As that pole moves to the left on the real axis, i.e., as a_0 becomes larger, the system becomes faster. This can be seen from the fact that the steady state value of the system response

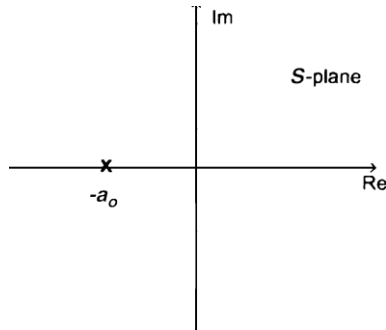
$$y_{ss} = \lim_{t \rightarrow \infty} y(t) = 1$$

is approached by the trajectory of $y(t)$ faster, as a_0 becomes larger. To see this, note that the value $1 - e^{-1}$ is attained at time $\tau = 1/a_0$, which is smaller as a_0 becomes larger. τ is the *time constant* of this first-order system; see below for further discussion of the time constant of a system.

We now derive the response of a second-order system to a unit step input (Fig. 19.1.7). Consider a system described by $\ddot{y}(t) + a_1 \dot{y}(t) + a_0 y(t) = a_0 u(t)$, which gives rise to the transfer function:

$$G(s) = \frac{a_0}{s^2 + a_1 s + a_0}$$

FIGURE 19.1.6 Pole location of a first-order system.



Notice that the steady-state value of the response to a unit step is

$$y_{ss} = \lim_{s \rightarrow 0} sG(s) = \frac{1}{1} = 1$$

note that this normalization or scaling to 1 is in fact the reason for selecting the constant numerator to be a_0 . $G(s)$ above does not have any finite zeros—only poles—as we want to study first the effect of the poles on the system behavior. We shall discuss the effect of adding a zero or an extra pole later.

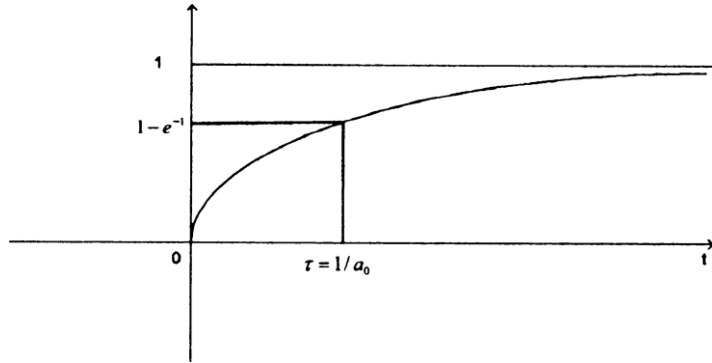


FIGURE 19.1.7 Step response of a first-order plant.

It is customary, and useful as we will see, to write the above transfer function as

$$G(s) = \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2}$$

where ζ is the *damping ratio* of the system and ω_n is the (*undamped*) *natural frequency* of the system, i.e., the frequency of oscillations when the damping is zero.

The poles of the system are

$$p_{1,2} = -\zeta\omega_n \pm \omega_n\sqrt{\zeta^2 - 1}$$

When $\zeta > 1$ the poles are real and distinct and the unit step response approaches its steady-state value of 1 without overshoot. In this case the system is *overdamped*. The system is called *critically damped* when $\zeta = 1$ in which case the poles are real, repeated, and located at $-\zeta\omega_n$.

The more interesting case is when the system is *underdamped* ($\zeta < 1$). In this case the poles are complex conjugate and are given by

$$p_{1,2} = -\zeta\omega_n \pm j\omega_n\sqrt{1 - \zeta^2} = \sigma + j\omega_d$$

The response to a unit step input in this case is

$$y(t) = \left[1 - \frac{e^{-\zeta\omega_n t}}{\sqrt{1 - \zeta^2}} \sin(\omega_d t + \theta) \right] q(t)$$

where $\theta = \cos^{-1} \zeta = \tan^{-1} (\sqrt{1 - \zeta^2} / \zeta)$, $\omega_d = \omega_n \sqrt{1 - \zeta^2}$, and $q(t)$ is the step function. The response to an *impulse input* ($u(t) = \delta(t)$) also called the *impulse response* $h(t)$ of the system is given in this case by

$$h(t) = \left[\frac{e^{-\zeta\omega_n t}}{\omega_n \sqrt{1 - \zeta^2}} \sin\left(\omega_n \sqrt{1 - \zeta^2} t\right) \right] q(t)$$

The second-order system is parameterized by the two parameters ζ and ω_n . Different choices for ζ and ω_n lead to different pole locations and to different behavior of the (modes of) the system. Fig. 19.1.8 shows the relation between the parameters and the pole location.

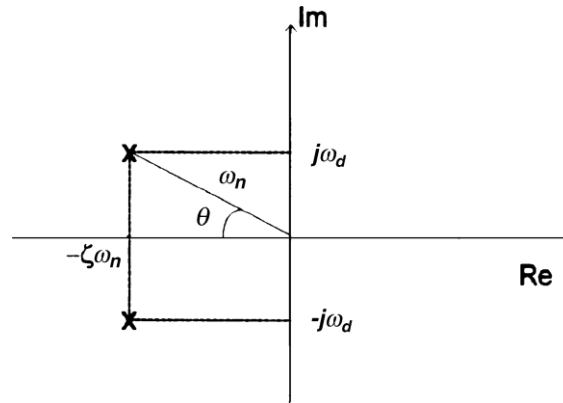


FIGURE 19.1.8 Relation between pole location and parameters.

Time Constant of a Mode and of a System. The time constant of a mode ce^{pt} of a system is the time value that makes $|pt| = 1$, i.e., $\tau = 1/|p|$. For example, in the above first-order system we have seen that $\tau = 1/a_0 = RC$. A pair of complex conjugate poles $p_{1,2} = \sigma \pm j\omega$ give rise to the term of the form $Ce^{\sigma t} \sin(\omega t + \theta)$. In this case, $\tau = 1/|\sigma|$, i.e., τ is again the inverse of the distance of the pole from the imaginary axis. The time constant of a system is the time constant of its dominant modes.

Transient Response Performance Specifications for a Second-Order Underdamped System

For the system

$$G(s) = \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2}$$

and a unit step input, explicit formulas for important measures of performance of its transient response can be derived. Note that the steady state is

$$y_{ss} = \lim_{s \rightarrow 0} sG(s) \frac{1}{s} = 1$$

The *rise time* t_r shows how long it takes for the system's output to rise from 0 to 66 percent of its final value (equal to 1 here) and it can be shown to be $t_r = (\pi - \theta)/\omega_n$, where $\theta = \cos^{-1} \zeta$ and $\omega_d = \omega_n \sqrt{1 - \zeta^2}$. The *settling time* t_s is the time required for the output to settle within some percentage, typically 2 percent or 5 percent, of its final value. $t_s \cong 4/\zeta\omega_n$ is the 2 percent settling time ($t_s \cong 3/\zeta\omega_n$ is the 5 percent settling time). Before the underdamped system settles, it will overshoot its final value. The *peak time* t_p measures the time it takes for the output to reach its first (and highest) peak value. M_p measures the actual *overshoot* that occurs in percentage terms of the final value. M_p occurs at time t_p , which is the time of the first and largest overshoot.

$$t_p = \frac{\pi}{\omega_d}, \quad M_p = 100e^{-\zeta\pi/\sqrt{1-\zeta^2}}\%$$

It is important to notice that the overshoot depends only on ζ . Typically, tolerable overshoot values are between 2.5 percent and 25 percent, which correspond to damping ratio ζ between 0.8 and 0.4.

Effect of Additional Poles and Zeros

The addition of an extra pole in the left-half s -plane (LHP) tends to slow the system down—the rise time of the system, for example, will become larger. When the pole is far to the left of the imaginary axis, its effect tends to be small. The effect becomes more pronounced as the pole moves toward the imaginary axis.

The addition of a zero in the LHP has the opposite effect, as it tends to speed the system up. Again the effect of a zero far away to the left of the imaginary axis tends to be small. It becomes more pronounced as the zero moves closer to the imaginary axis.

The addition of a zero in the right-half s -plane (RHP) has a delaying effect much more severe than the addition of a LHP pole. In fact a RHP zero causes the response (say, to a step input) to start toward the wrong direction. It will move down first and become negative, for example, before it becomes positive again and starts toward its steady-state value. Systems with RHP zeros are called *nonminimum phase systems* (for reasons that will become clearer after the discussion of the frequency design methods) and are typically difficult to control. Systems with only LHP poles (stable) and LHP zeros are called *minimum phase systems*.

CLASSICAL CONTROL DESIGN METHODS

In this section, we focus on the problem of controlling a single-input and single-output (SISO) LTI plant. It is understood from the above sections that such a plant can be represented by a transfer function $G_p(s)$. The closed-loop system is shown in Fig. 19.1.9.

The goal of feedback control is to make the output of the plant, y , follow the reference input r as closely as possible. Classical design methods are those used to determine the controller transfer function $G_c(s)$ so that the closed-loop system, represented by the transfer function:

$$G_{CL}(s) = \frac{G_c(s)G_p(s)}{1 + G_c(s)G_p(s)}$$

has desired characteristics.

Design Specifications and Constraints

The design specifications are typically described in terms of step response, i.e., r is the set point described as a step-like function. These specifications are given in terms of transient response and steady-state error, assuming the feedback control system is stable. The transient response is characterized by the rise time, i.e., the time it takes for the output to reach 66 percent of its final value, the settling time, i.e., the time it takes for the output to settle within 2 percent of its final value, and the percent overshoot, which is how much the output exceeds the set-point r percentagewise during the period that y converges to r . The steady-state error refers to the difference, if any, between y and r as y reaches its steady-state value.

There are many constraints a control designer has to deal with in practice, as shown in Fig. 19.1.10. They can be described as follows:

1. **Actuator Saturation:** The input u to the plant is physically limited to a certain range, beyond which it “saturates,” i.e., becomes a constant.
2. **Disturbance Rejection and Sensor Noise Reduction:** There are always disturbances and sensor noises in the plant to be dealt with.

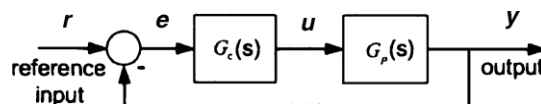


FIGURE 19.1.9 Feedback control configuration.

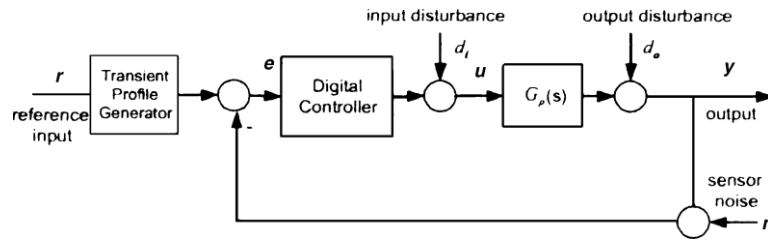


FIGURE 19.1.10 Closed-loop simulator setup.

3. *Dynamic Changes in the Plant:* Physical systems are almost never truly linear nor time invariant.
4. *Transient Profile:* In practice, it is often not enough to just move y from one operating point to another. How it gets there is sometimes just as important. Transient profile is a mechanism to define the desired trajectory of y in transition, which is of great practical concerns. The smoothness of y and its derivatives, the energy consumed, the maximum value and the rate of change required of the control action are all influenced by the choice of transient profile.
5. *Digital Control:* Most controllers are implemented today in digital forms, which makes the sampling rate and quantization errors limiting factors in the controller performance.

Control Design Strategy Overview

The control strategies are summarized here in ascending order of complexity and, hopefully, performance.

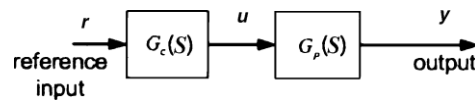


FIGURE 19.1.11 Open-loop control configuration.

1. *Open-Loop Control:* If the plant transfer function is known and there is very little disturbance, a simple open loop controller, as shown in Fig. 19.1.11, would satisfy most design requirements.

Where $G_c(s)$ is an approximate inverse of $G_p(s)$. Such control strategy has been used as an economic means in controlling stepper motors, for example.

2. *Feedback Control with a Constant Gain:* With significant disturbance and dynamic variations in the plant, feedback control, as shown in Fig. 19.1.9 is the only choice; see also Appendix B. Its simplest form is $G_c(s) = k$, or $u = ke$, where k is a constant. Such proportional controller is very appealing because of its simplicity. The common problems with this controller are significant steady-state error and overshoot.
3. *Proportional-Integral-Derivative Controller:* To correct the above problems with the constant gain controller, two additional terms are added:

$$u = k_p e + k_i \int e + k_d \dot{e} \quad \text{or} \quad G_c(s) = k_p + k_i / s + k_d s$$

This is the well-known PID controller, which is used by most engineers in industry today. The design can be quite intuitive: the proportional term usually plays the key role, with the integral term added to reduce/eliminate the steady-state error and the derivative term the overshoot. The primary drawbacks of PID is that the integrator introduces phase lag that could lead to stability problems and the differentiator makes the controller sensitive to noise.

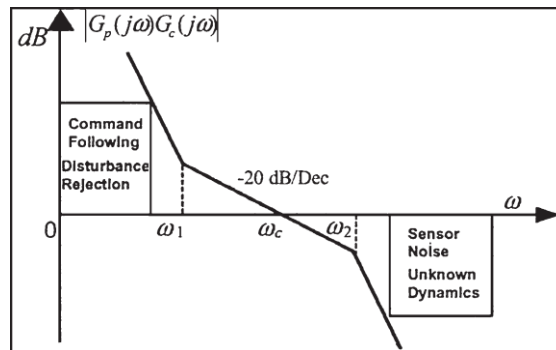


FIGURE 19.1.12 Loop-shaping.

4. **Root Locus Method:** A significant portion of most current control textbooks is devoted to the question of how to place the poles of the closed-loop system in Fig. 19.1.9 at desired locations, assuming we know where they are. Root locus is a graphical technique to manipulate the closed-loop poles given the open-loop transfer function. This technique is most effective if disturbance rejection, plant dynamical variations, and sensor noise are not to be considered. This is because these properties cannot be easily linked to closed loop pole locations.
5. **Loop-Shaping Method:** Loop-shaping [5] refers to the manipulation of the *loop gain* frequency response, $L(j\omega) = G_p(j\omega)G_c(j\omega)$, as a control design tool. It is the only existing design method that can bring most of design specifications and constraints, as discussed above, under one umbrella and systematically find a solution. This makes it a very useful tool in understanding, diagnosing, and solving practical control problems. The loop-shaping process consists of two steps:
 - a. Convert all design specifications to loop gain constraints, as shown in Fig. 19.1.12.
 - b. Find a controller $G_c(s)$ to meet the specifications.

Loop-shaping as a concept and a design tool helped the practicing engineers greatly in improving the PID loop performance and stability margins. For example, a PID implemented as a lead-lag compensator is commonly seen in industry today. This is where the classical control theory provides the mathematical and design insights on why and how feedback control works. It has also laid the foundation for modern control theory.

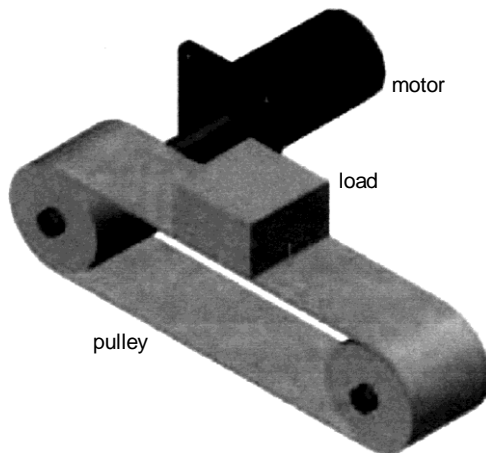


FIGURE 19.1.13 A Digital servo control design example.

Example Consider a motion control system as shown in Fig. 19.1.13 below. It consists of a digital controller, a dc motor drive (motor and power amplifier), and a load of 235 lb that is to be moved linearly by 12 in. in 0.3 s with an accuracy of 1 percent or better. A belt and pulley mechanism is used to convert the motor rotation a linear motion. Here a servo motor is used to drive the load to perform a linear motion. The motor is coupled with the load through a pulley.

The design process involves:

1. Selection of components including motor, power amplifier, the belt-and-pulley, the feedback devices (position sensor and/or speed sensor)
2. Modeling of the plant
3. Control design and simulation
4. Implementation and tuning

The first step results in a system with the following parameters:

1. *Electrical:*

- Winding resistance and inductance: $R_a = 0.4$ mho $L_a = 8$ mH (the transfer function of armature voltage to current is $(1/R_a)/[(L_a/R_a)s + 1]$)
- back emf constant: $K_E = 1.49$ V/(rad/s)
- power amplifier gain: $K_{pa} = 80$
- current feedback gain: $K_{cf} = 0.075$ V/A

2. *Mechanical:*

- Torque constant: $K_t = 13.2$ in-lb/A
- Motor inertia $J_m = .05$ lb-in.s²
- Pulley radius $R_p = 1.25$ in.
- Load weight: $W = 235$ lb (including the assembly)
- Total inertia $J_t = J_m + J_l = 0.05 + (W/g)R_p^2 = 1.0$ lb-in.s²

With the maximum armature current set at 100 A, the maximum Torque $= K_t I_{a,\max} = 13.2 \times 100 = 1320$ in.-lb; the maximum angular acceleration $= 1320/J_t = 1320$ rad/s², and the maximum linear acceleration $= 1320 \times R_p = 1650$ in./s² $= 4.27$ g's (1650/386). As it turned out, they are sufficient for this application.

The second step produces a simulation model (Fig. 19.1.14).

A simplified transfer function of the plant, from the control input, v_c (in volts), to the linear position output, x_{out} (in inches), is

$$G_p(s) = \frac{206}{s(s+3)}$$

An open loop controller is not suitable here because it cannot handle the torque disturbances and the inertia change in the load. Now consider the feedback control scheme in Fig. 19.1.9 with a constant controller, $u = ke$. The root locus plot in Fig. 19.1.15 indicates that, even at a high gain, the real part of the closed-loop poles does not exceed -1.5 , which corresponds to a settling time of about 2.7 s. This is far slower than desired.

In order to make the system respond faster, the closed-loop poles must be moved further away from the $j\omega$ axis. In particular, a settling time of 0.3 s or less corresponds to the closed-loop poles with real parts smaller than -13.3 . This is achieved by using a PD controller of the form

$$G_c(s) = K(s+3); K \geq 13.3/206$$

will result in a settling time of less than 0.3 s.

The above PD design is a simple solution in servo design that is commonly used. There are several issues, however, that cannot be completely resolved in this framework:

1. Low-frequency torque disturbance induces steady-state error that affects the accuracy.
2. The presence of a resonant mode within or close to the bandwidth of the servo loop may create undesirable vibrations.
3. Sensor noise may cause the control signal to be very noisy.
4. The change in the dynamics of the plant, for example, the inertia of the load, may require frequency tweaking of the controller parameters.
5. The step-like set point change results in an initial surge in the control signal and could shorten the life span of the motor and other mechanical parts.

These are problems that most control textbook do not adequately address, but they are of significant importance in practice. The first three problems can be tackled using the loop-shaping design technique introduced above. The tuning problem is an industrywide design issue and the focus of various research and development efforts. The last problem is addressed by employing a smooth transient as the set point, instead of a step-like set-point. This is known as the "motion profile" in industry.

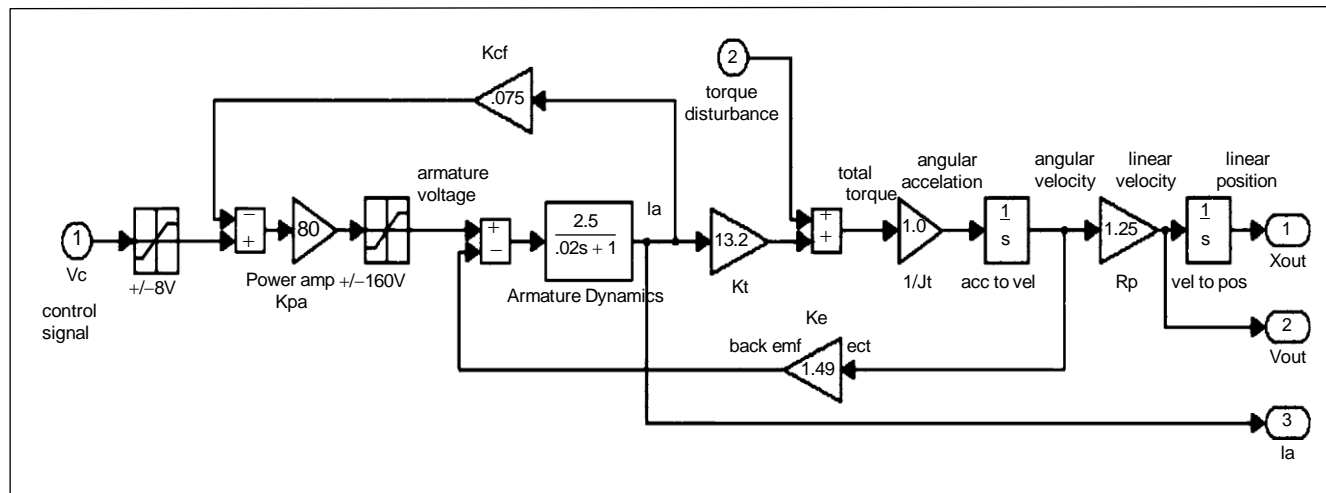


FIGURE 19.1.14 Simulation model of the motion control system.

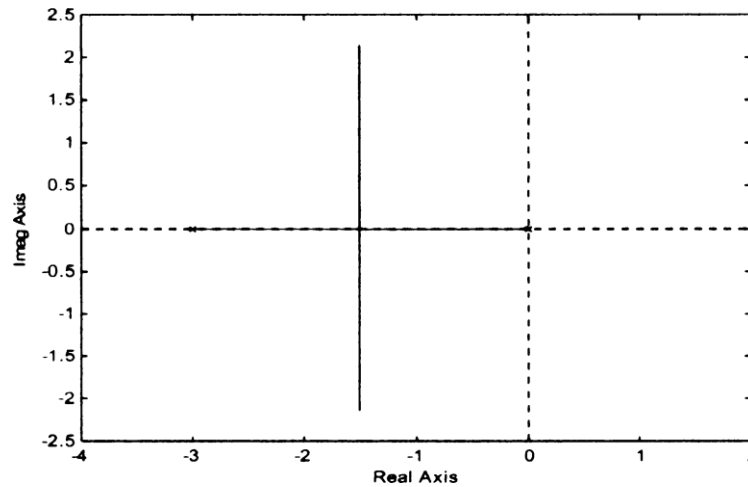


FIGURE 19.1.15 Root locus plot of the servo design problem.

Evaluation of Control Systems

Analysis of control system provides crucial insights to control practitioners on why and how feedback control works. Although the use of PID precedes the birth of classical control theory of the 1950s by at least two decades, it is the latter that established the control engineering discipline. The core of classical control theory are the frequency-response-based analysis techniques, namely, Bode and Nyquist plots, stability margins, and so forth.

In particular, by examining the loop gain frequency response of the system in Fig. 19.1.9, that is, $L(j\omega) = G_c(j\omega)G_p(j\omega)$, and the sensitivity function $1/[1 + L(j\omega)]$, one can determine the following:

1. How fast the control system responds to the command or disturbance input (i.e., the bandwidth).
2. Whether the closed-loop system is stable (Nyquist Stability Theorem); If it is stable, how much dynamic variation it takes to make the system unstable (in terms of the gain and phase change in the plant). It leads to the definition of gain and phase margins. More broadly, it defines how robust the control system is.
3. How sensitive the performance (or closed-loop transfer function) is to the changes in the parameters of the plant transfer function (described by the sensitivity function).
4. The frequency range and the amount of attenuation for the input and output disturbances shown in Fig. 19.1.10 (again described by the sensitivity function).

Evidently, these characteristics obtained via frequency-response analysis are invaluable to control engineers. The efforts to improve these characteristics led to the lead-lag compensator design and, eventually, loop-shaping technique described above.

Example: The PD controller in Fig. 19.1.10 is known to be sensitive to sensor noises. A practical cure to this problem is: add a low pass filter to the controller to attenuate high-frequency noises, that is,

$$G_c(s) = \frac{13.3(s+3)}{206 \left(\frac{s}{133} + 1 \right)^2}$$

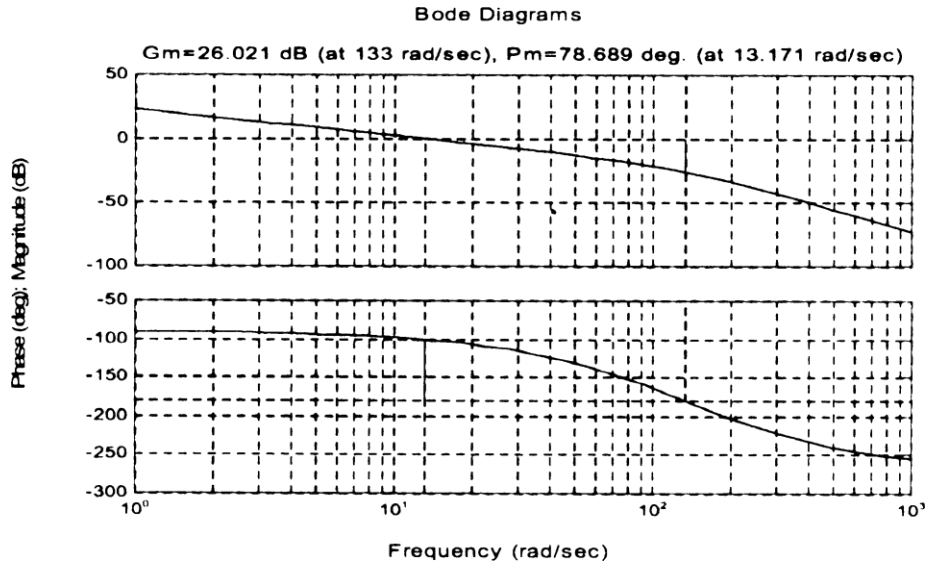


FIGURE 19.1.16 Bode plot evaluation of the control design.

The loop gain transfer function is now

$$L(s) = G_p(s)G_c(s) = \left(\frac{13.3}{s} \right)^2 \frac{1}{s(133s + 1)}$$

The bandwidth of the low pass filter is chosen to be one decade higher than the loop gain bandwidth to maintain proper gain and phase margins. The Bode plot of the new loop gain, as shown in Fig. 19.1.16, indicates that (a) the feedback system has a bandwidth 13.2 rad/s, which corresponds to a 0.3 s settling time as specified and (b) this design has adequate stability margins (gain margin is 26 dB and phase margin is 79°).

Digital Implementation

Once the controller is designed and simulated successfully, the next step is to digitize it so that it can be programmed into the processor in the digital control hardware. To do this:

1. Determine the sampling period T_s and the number of bits used in analog-to-digital converter (ADC) and digital-to-analog converter (DAC).
2. Convert the continuous time transfer function $G_c(s)$ to its corresponding form in discrete time transfer function $G_{cd}(z)$ using, for example, the Tustin's method, $s = (1/T_s)(z - 1)/(z + 1)$.
3. From $G_{cd}(z)$, derive the difference equation, $u(k) = g(u(k-1), u(k-2), \dots, y(k), y(k-1), \dots)$, where g is a linear algebraic function.

After the conversion, the sampled data system, with the plant running in continuous time and the controller in discrete time, should be verified in simulation first before the actual implementation. The quantization error and sensor noise should also be included to make it realistic.

The minimum sampling frequency required for a given control system design has not been established analytically. The rule of thumb given in control textbooks is that $f_s = 1/T_s$ should be chosen approximately 30 to 60 times the bandwidth of the closed-loop system. Lower-sampling frequency is possible after careful tuning but the aliasing, or signal distortion, will occur when the data to be sampled have significant energy above the

Nyquist frequency. For this reason, an antialiasing filter is often placed in front of the ADC to filter out the high-frequency contents in the signal.

Typical ADC and DAC chips have 8, 12, and 16 bits of resolution. It is the length of the binary number used to approximate an analog one. The selection of the resolution depends on the noise level in the sensor signal and the accuracy specification. For example, the sensor noise level, say 0.1 percent, must be below the accuracy specification, say 0.5 percent. Allowing one bit for the sign, an 8-bit ADC with a resolution of $1/2^7$, or 0.8 percent, is not good enough; similarly, a 16-bit ADC with a resolution of 0.003 percent is unnecessary because several bits are “lost” in the sensor noise. Therefore, a 12-bit ADC, which has a resolution of 0.04 percent, is appropriate for this case. This is an example of “error budget,” as it is known among designers, where components are selected economically so that the sources of inaccuracies are distributed evenly.

Converting $G_c(s)$ to $G_{cd}(z)$ is a matter of numerical integration. There have been many methods suggested, some are too simple and inaccurate (such as the Euler’s forward and backward methods), others are too complex. The Tustin’s method suggested above, also known as trapezoidal method or bilinear transformation, is a good compromise. Once the discrete transfer function $G_{cd}(z)$ is obtained, finding the corresponding difference equation that can be easily programmed in C is straightforward. For example, given a controller with input $e(k)$ and output $u(k)$, and the transfer function

$$G_{cd}(z) = \frac{z+2}{z+1} = \frac{1+2z^{-1}}{1+z^{-1}}$$

the corresponding input-output relationship is

$$u(k) = \frac{1+2z^{-1}}{1+z^{-1}} e(k)$$

or equivalently, $(1+z^{-1})u(k) = (1+2z^{-1})e(k)$. That is, $u(k) = -u(k-1) + e(k) + 2e(k-1)$.

Finally, the presence of the sensor noise usually requires that an antialiasing filter be used in front of the ADC to avoid distortion of the signal in ADC. The phase lag from such a filter must not occur at the crossover frequency (bandwidth) or it will reduce the stability margin or even destabilize the system. This puts yet another constraint on the controller design.

ALTERNATIVE DESIGN METHODS

Nonlinear PID

Using nonlinear PID (NPID) is an alternative to PID for better performance. It maintains the simplicity and intuition of PID, but empowers it with nonlinear gains. An example of NPID is shown in Fig. 19.1.17. The need for the integral control is reduced, by making the proportional gain larger, when the error is small. The limited

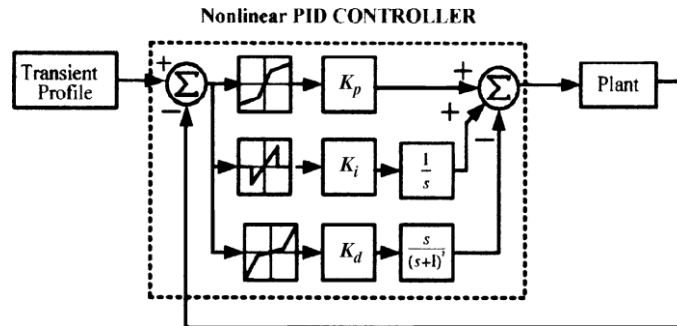


FIGURE 19.1.17 Nonlinear PID for a power converter control problem.

authority integral control has its gain zeroed outside a small interval around the origin to reduce the phase lag. Finally the differential gain is reduced for small errors to reduce sensitivities to sensor noise. See Ref. 8.

State Feedback and Observer-Based Design

If the state space model of the plant

$$\begin{aligned}\dot{x} &= Ax + Bu \\ y &= Cx + Du\end{aligned}$$

is available, the pole-placement design can be achieved via state feedback

$$u = r + Kx$$

where K is the gain vector to be determined so that the eigenvalues of the closed-loops system

$$\begin{aligned}\dot{x} &= (A + BK)x + Br \\ y &= Cx + Du\end{aligned}$$

are at the desired locations, assuming they are known. Usually the state vector is not available through measurements and the state observer is of the form

$$\begin{aligned}\dot{\hat{x}} &= A\hat{x} + Bu + L(y - \hat{y}) \\ \hat{y} &= C\hat{x} + Du\end{aligned}$$

where \hat{x} is the estimate of x and L is the observer gain vector to be determined.

The state feedback design approach has the same drawbacks as those of Root Locus approach, but the use of the state observer does provide a means to extract the information about the plant that is otherwise unavailable in the previous control schemes, which are based on the input-output descriptions of the plant. This proves to be valuable in many applications. In addition, the state space methodologies are also applicable to systems with many inputs and outputs.

Controllability and Observability. Controllability and observability are useful system properties and are defined as follows. Consider an n th order system described by

$$\dot{x} = Ax + Bu, \quad z = Mx$$

where A is an $n \times n$ matrix. The system is *controllable* if it is possible to transfer the state to any other state in finite time. This property is important as it measures, for example, the ability of a satellite system to reorient itself to face another part of the earth's surface using the available thrusters; or to shift the temperature in an industrial oven to a specified temperature. Two equivalent tests for controllability are:

The system (or the pair (A, B)) is *controllable* if and only if the controllability matrix $C = [B, AB, \dots, A^{n-1}B]$ has full (row) rank n . Equivalently if and only if $[s_i I - A, B]$ has full (row) rank n for all eigenvalues s_i of A . The system is *observable* if by observing the output and the input over a finite period of time it is possible to deduce the value of the state vector of the system. If, for example, a circuit is observable it may be possible to determine all the voltages across the capacitors and all currents through the inductances by observing the input and output voltages.

The system (or the pair (A, C)) is *observable* if and only if the observability matrix

$$\theta = \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{n-1} \end{bmatrix}$$

has full (column) rank n . Equivalently if and only if

$$\begin{bmatrix} s_i I - A \\ C \end{bmatrix}$$

has full (column) rank n for all eigenvalues s_i of A .

Consider now the transfer function

$$G(s) = C(sI - A)^{-1}B + D$$

Note that, by definition, in a transfer function all possible cancellations between numerator and denominator polynomials are assumed to have already taken place. In general, therefore, the poles of $G(s)$ are some (or all) of the eigenvalues of A . It can be shown that when the system is both controllable and observable no cancellations take place and so in this case the poles of $G(s)$ are exactly the eigenvalues of A .

Eigenvalue Assignment Design. Consider the equations: $\dot{x} = Ax + Bu$, $y = Cx + Du$, and $u = p + kx$. When the system is controllable, K can be selected to assign the closed-loop eigenvalues to any desired locations (real or complex conjugate) and thus significantly modify the behavior of the open-loop system. Many algorithms exist to determine such K . In the case of a single input, there is a convenient formula called Ackermann's formula

$$K = -[0, \dots, 0, 1] C^{-1} \alpha_d(A)$$

where $C = [B, \dots, A^{n-1}B]$ is the $n \times n$ controllability matrix and the roots of $\alpha_d(s)$ are the desired closed-loop eigenvalues.

Example Let

$$A = \begin{bmatrix} 1/2 & 1 \\ 1 & 2 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

and the desired eigenvalues be $-1 \pm j$

Here

$$C = [B, AB] = \begin{bmatrix} 1 & 3/2 \\ 1 & 3 \end{bmatrix}$$

Note that A has eigenvalues at 0 and 5/2. We wish to determine K so that the eigenvalues of $A + BK$ are at $-1 \pm j$, which are the roots of $\alpha_d(s) = s^2 + 2s + 2$.

Here

$$C = [B, AB] = \begin{bmatrix} 1 & 3/2 \\ 1 & 3 \end{bmatrix}$$

and

$$\alpha_d(A) = A^2 + 2A + 2I = \left(\begin{bmatrix} 1/2 & 1 \\ 1 & 2 \end{bmatrix}^2 + 2 \begin{bmatrix} 1/2 & 1 \\ 1 & 2 \end{bmatrix} + 2 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) = \begin{bmatrix} 17/4 & 9/2 \\ 9/2 & 11 \end{bmatrix}$$

Then

$$K = -[0 \ 1]C^{-1}\alpha_d(A) = [-1/6 \ -13/3]$$

Here

$$A + BK = \begin{bmatrix} 1/3 & -10/3 \\ 5/6 & -7/3 \end{bmatrix}$$

which has the desired eigenvalues.

Linear Quadratic Regulator (LQR) Problem. Consider

$$\dot{x} = Ax + Bu, \quad z = Mx$$

We wish to determine $u(t)$, $t \geq 0$, which minimizes the quadratic cost

$$J(u) = \int_0^\infty \left[x^T(t) (M^T Q M) x + u^T(t) R u(t) \right] dt$$

for any initial state $x(0)$. The weighting matrices Q and R are real, symmetric ($Q = Q^T$, $R = R^T$), Q and R are positive definite ($R > 0$, $Q > 0$) and $M^T Q M$ is positive semidefinite ($M^T Q M \geq 0$). Since $R > 0$, the term $u^T R u$ is always positive for any $u \neq 0$, by definition. Minimizing its integral forces $u(t)$ to remain small. $M^T Q M \geq 0$ implies that $x^T M^T Q M x$ is positive, but it can also be zero for some $x \neq 0$, this allows some of the states to be treated as “do not care states.” Minimizing the integral of $x^T M^T Q M x$ forces the states to become smaller as time progresses. It is convenient to take Q (and R in the multi-input case) to be diagonal with positive entries on the diagonal. The above performance index is designed so that the minimizing control input drives the states to the zero state, or as close as possible, without using excessive control action, in fact minimizing the control energy. When $(A, B, Q^{1/2}M)$ is controllable and observable, the solution $u^*(t)$ of this optimal control problem is a state feedback control law, namely,

$$u^*(t) = K^* x(t) = -R^{-1} B^T P_c^* x(t)$$

where P_c^* is the unique symmetric positive definite solution of the *algebraic Riccati equation*:

$$A^T P_c + P_c A - P_c B R^{-1} B^T P_c + M^T Q M = 0$$

Example. Consider

$$\dot{x} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix} x + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u, \quad y = [1 \ 0] x$$

And let

$$J = \int_0^\infty \left(y^2(t) + 4u^2(t) \right) dt$$

Here

$$M = C, \quad Q = 1, \quad M^T Q M = C^T C = \begin{bmatrix} 1 \\ 0 \end{bmatrix} [1 \ 0] = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad R = 4$$

Solving the Riccati equation we obtain

$$P^* = \begin{bmatrix} 2 & 2\sqrt{2} \\ 2\sqrt{2} & 2 \end{bmatrix}$$

and

$$u^*(t) = K^* x(t) = -\frac{1}{4} \begin{bmatrix} 2 & 2 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 2\sqrt{2} \\ 2\sqrt{2} & 2 \end{bmatrix} x(t) = -\frac{1}{2} \begin{bmatrix} 1 & \sqrt{2} \\ \sqrt{2} & 1 \end{bmatrix} x(t)$$

Linear State Observers. Since the states of a system contain a great deal of useful information, knowledge of the state vector is desirable. Frequently, however, it may be either impossible or impractical to obtain measurements of all states. Therefore, it is important to be able to estimate the states from available measurements, namely, of inputs and outputs.

Let the system be

$$\dot{x} = Ax + Bu, \quad y = Cx + Du$$

An asymptotic state estimator of the full state, also called *Luenberger observer*, is given by

$$\dot{\hat{x}} = A\hat{x} + Bu + L(y - \hat{y})$$

where L is selected so that all eigenvalues of $A - LC$ are in the LHP (have negative real parts). Note that a L that arbitrarily assigns the eigenvalues of $A - LC$ exists if and only if the system is observable. The observer may be written as

$$\hat{x} = (A - LC)\hat{x} + [B - LD, K] \begin{bmatrix} u \\ y \end{bmatrix}$$

which clearly shows the role of u and y ; they are the inputs to the observer. If the error $e(t) = x(t) - \hat{x}(t)$ then $e(t) = e^{(A-LC)t}e(0)$, which shows that $e(t) \rightarrow 0$ or $\hat{x}(t) \rightarrow x(t)$ as $t \rightarrow \infty$.

To determine appropriate L , note that $(A - LC)^T = A^T + C^T(-L) = A + B\bar{K}$, which is the problem addressed above in the state feedback case. One could also use the following observable version of Ackermann's formula, namely,

$$L = \alpha_d(A) \theta^{-1} [0, \dots, 0, 1]^T$$

where

$$\theta = \begin{bmatrix} C \\ CA \\ \vdots \\ CA^{n-1} \end{bmatrix}$$

The gain L in the above estimator may be determined so that it is optimal in an appropriate sense. In the following, some of the key equations of such an optimal estimator (*Linear Quadratic Gaussian* (LQG)), also known as the *Kalman-Bucy filter*, are briefly outlined.

Consider

$$\dot{x} = Ax + Bu + \Gamma w, \quad y = Cx + v$$

where w and v represent process and measurement noise terms. Both w and v are assumed to be white, zero-mean Gaussian stochastic processes, i.e., they are uncorrelated in time and have expected values $E[w] = 0$ and $E[v] = 0$. Let $E[ww^T] = W$, $E[vv^T] = V$ denote the covariances where W , V are real, symmetric and positive definite matrices. Assume also that the noise processes w and v are independent, i.e., $E[wv^T] = 0$. Also assume that the initial state $x(0)$ is a Gaussian random variable of known mean, $E[x(0)] = x_0$, and known covariance $E[(x(0) - x_0)(x(0) - x_0)^T] = P_{e0}$. Assume also that $x(0)$ is independent of w and v .

Consider now the estimator

$$\dot{\hat{x}} = (A - LC)\hat{x} + Bu + Ly$$

and let $(A, \Gamma W^{1/2}, C)$ be controllable and observable. Then the error covariance $E[(x - \hat{x})(x - \hat{x})^T]$ is minimized when the filter gain $L = P_e^* C^T V^{-1}$, where P_e^* denotes the symmetric, positive definite solution of the (dual to control) algebraic Riccati equation

$$P_e A^T + A P_e - P_e C^T V^{-1} C P_e + \Gamma W \Gamma^T = 0$$

The above Riccati is the *dual* to the Riccati equation for optimal control and can be obtained from the optimal control equation by making use of the substitutions:

$$A \rightarrow A^T, B \rightarrow C^T, M \rightarrow \Gamma^T, R \rightarrow V, Q \rightarrow W$$

In the state feedback control law $u = Kx + r$, when state measurements are not available, it is common to use the estimate of state \hat{x} from a Luenberger observer. That is, given

$$x = Ax + Bu, \quad y = Cx + Du$$

the control law is $u = K\hat{x} + r$, where \hat{x} is the state estimate from the observer

$$\dot{\hat{x}} = (A - LC)\hat{x} + [B - KD, K] \begin{bmatrix} u \\ y \end{bmatrix}$$

The closed-loop system is then of order $2n$ since the plant and the observer are each of order n . It can be shown that in this case, of linear output feedback control design, the design of the control law and of the gain K (using, for example, LQR) can be carried out independent of the design of the estimator and the filter gain L (using, for example, LQG). This is known as the *separation property*.

It is remarkable to notice that the overall transfer function of the compensated system that includes the state feedback and the observer is

$$T(s) = (C + DK)[sI - (A + BK)]^{-1} B + D$$

which is exactly the transfer function one would obtain if the state x were measured directly and the state observer were not present. This is of course assuming zero initial conditions (to obtain the transfer function); if nonzero initial conditions are present, then there is some deterioration of performance owing to observer dynamics, and the fact that at least initially the state estimate typically contains significant error.

ADVANCED ANALYSIS AND DESIGN TECHNIQUES

This section covered some fundamental analysis and design methods in classical control theory, the development of which was primarily driven by engineering practice and needs. Over the last few decades, vast efforts in control research have led to the creation of modern mathematical control theory, or advanced control, or control science. This development started with optimal control theory in the 50s and 60s to study the optimality of control design; a brief glimpse of optimal control was given above. In optimal control theory, a cost function is to be minimized, and analytical or computational methods are used to derive optimal controllers. Examples include minimum fuel problem, time-optimal control (Bang-Bang) problem, LQ, H_2 , and H_∞ , each corresponding to a different cost function. Other major branches in modern control theory include multi-input multi-output (MIMO) control systems methodologies, which attempt to extend well-known SISO design methods and concepts to MIMO problems; adaptive control, designed to extend the operating range of a controller by adjusting automatically the controller parameters based on the estimated dynamic changes in the plants; analysis and design of nonlinear control systems, and so forth.

A key problem is the robustness of the control system. The analysis and design methods in control theory are all based on the mathematical model of the plants, which is an approximate description of physical processes. Whether a control system can tolerate the uncertainties in the dynamics of the plants, or how much uncertainty it takes to make a system unstable, is studied in robust control theory, where H_2 , H_∞ , and other analysis and design methods were originated. Even with recent progress, open problems remain when dealing with real world applications. Some recent approaches such as in Ref. 8 attempt to address some of these difficulties in a realistic way.

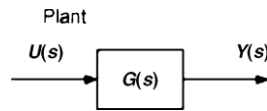
APPENDIX: OPEN AND CLOSED LOOP STABILIZATION

It is impossible to stabilize an unstable system using open-loop control, owing to system uncertainties. In general, closed loop or feedback control is necessary to control a system—stabilize if unstable and improve performance—because of uncertainties that are always present. Feedback provides current information about the system and so the controller does not have to rely solely on incomplete system information contained in a nominal plant model. These uncertainties are system parameter uncertainties and also uncertainties induced on the system by its environment, including uncertainties in the initial condition of the system, and uncertainties because of disturbances and noise.

Consider the plant with transfer function

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

where the pole location at 1 is inaccurately known.



The corresponding description in the time domain using differential equations is $\dot{y}(t) - (1 + \varepsilon)y(t) = u(t)$. Solving, using Laplace transform, we obtain $sY(s) - y(0) - (1 + \varepsilon)Y(s) = U(s)$ from which

$$Y(s) = \frac{y(0)}{s - (1 + \varepsilon)} + \frac{1}{s - (1 + \varepsilon)}U(s)$$

Consider now the controller with transfer function

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

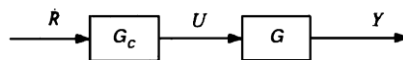


The corresponding description in the time domain using differential equations is $\dot{u}(t) + 2u(t) = \dot{r}(t) - r(t)$. Solving, using Laplace transform, we obtain $sU(s) - u(0) + 2U(s) = sR(s) - R(s)$ from which

$$U(s) = \frac{u(0)}{s + 2} + \frac{s - 1}{s + 2}R(s)$$

Connect now the plant and the controller in series (open-loop control)

Connecting in Series - Open loop



The overall transfer function is

$$T = GG_c = \frac{s-1}{[s-(1+\varepsilon)](s+2)}$$

Including the initial conditions

$$Y(s) = \frac{(s+2)y(0) + u(0)}{[s-(1+\varepsilon)](s+2)} + \frac{s-1}{[s-(1+\varepsilon)](s+2)} R(s)$$

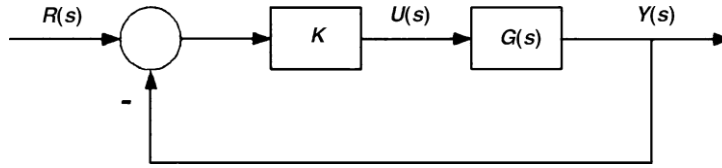
It is now clear that open-loop control cannot be used to stabilize the plant:

1. First, because of the uncertainties in the plant parameters. Note that the plant pole is not exactly at +1 but at $1 + \varepsilon$ and so the controller zero cannot cancel the plant pole exactly.
2. Second, even if we had knowledge of the exact pole location, that is, $\varepsilon = 0$, and

$$Y(s) = \frac{(s+2)y(0) + r(0)}{(s-1)(s+2)} + \frac{1}{s+2} R(s)$$

still we cannot stabilize the system because of the uncertainty in the initial conditions. We cannot, for example, select $r(0)$ so as to cancel the unstable pole at +1 because $y(0)$ may not be known exactly.

We shall now stabilize the above plant using a simple feedback controller.



Consider a unity feedback control system with the controller being just a gain k to be determined. The closed-loop transfer function is

$$T(s) = \frac{kG(s)}{1 + kG(s)} = \frac{k}{s - (1 + \varepsilon) + k}$$

Working in the time domain,

$$\dot{y} - (1 + \varepsilon)y = u = k(r - y)$$

from which

$$\dot{y} + [k - (1 + \varepsilon)]y = kr$$

Using Laplace transform we obtain

$$sY(s) - y(0) + [k - (1 + \varepsilon)]Y(s) = kR(s)$$

and

$$Y(s) = \frac{y(0)}{s + k - (1 + \varepsilon)} + \frac{k}{s + k - (1 + \varepsilon)} R(s)$$

It is now clear that if the controller gain is selected so that $k > 1 + \varepsilon$, then the system is stable. In fact, we could have worked with the nominal system to derive $k > 1$ for stability. For stability robustness, we take k somewhat larger than 1 to have some safety margin and satisfy $k > 1 + \varepsilon$ for some unknown small ε .

REFERENCES

1. Dorf, R. C., and R. H. Bishop, "Modern Control Systems," 9th ed., Prentice Hall, 2001.
2. Franklin, G. F., J. D. Powell, and A. Emami-Naeimi, "Feedback Control of Dynamic Systems," 3rd ed., Addison-Wesley, 1994.
3. Kuo, B. C., "Automatic Control Systems," 7th ed., Prentice Hall, 1995.
4. Ogata, K., "Modern Control Engineering," 3rd ed., Prentice Hall, 1997.
5. Rohrs, C. E., J. L. Melsa, and D. G. Schultz, "Linear Control Systems," McGraw-Hill, 1993.
6. Antsaklis, P. J., and A. N. Michel, "Linear Systems," McGraw-Hill, 1997.
7. Goodwin, G. C., S. F. Graebe, and M. E. Salgado, "Control System Design," Prentice Hall, 2001.
8. Gao, Z., Y. Huang, and J. Han, "An Alternative Paradigm for Control System Design," Presented at the 40th IEEE Conference on Decision and Control, Dec 4–7, 2001, Orlando, FL.

