

Fig. 8.3 Discrete-time impulse function.

1. Discrete-Time Impulse Function $\delta[k]$

The discrete-time counterpart of the continuous-time impulse function $\delta(t)$ is $\delta[k]$, defined by

$$\delta[k] = \begin{cases} 1 & k = 0 \\ 0 & k \neq 0 \end{cases} \quad (8.1)$$

This function, also called the unit impulse sequence, is shown in Fig. 8.3a. The time-shifted impulse sequence $\delta[k-m]$ is depicted in Fig. 8.3b. Unlike its continuous-time counterpart $\delta(t)$, this is a very simple function without any mystery.

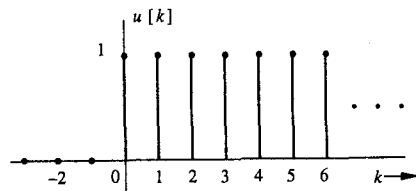
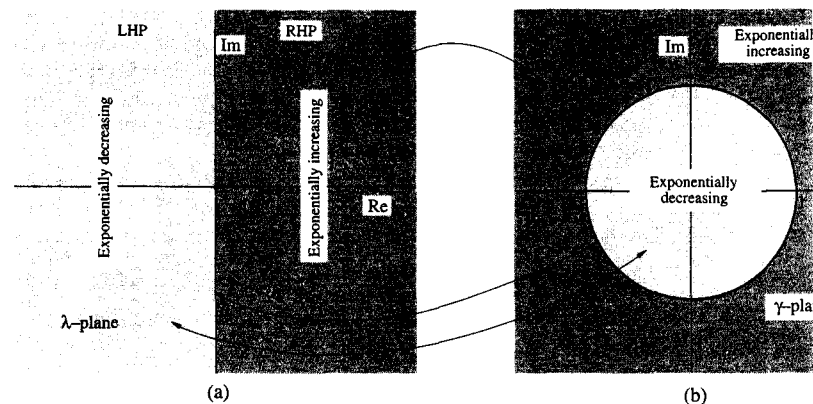
Later, we shall express an arbitrary input $f[k]$ in terms of impulse components. The (zero-state) system response to input $f[k]$ can then be obtained as the sum of system responses to impulse components of $f[k]$.

2. Discrete-Time Unit Step Function $u[k]$

The discrete-time counterpart of the unit step function $u(t)$ is $u[k]$ (Fig. 8.4), defined by

$$u[k] = \begin{cases} 1 & \text{for } k \geq 0 \\ 0 & \text{for } k < 0 \end{cases} \quad (8.2)$$

If we want a signal to start at $k = 0$ (so that it has a zero value for all $k < 0$), we need only multiply the signal with $u[k]$.

Fig. 8.4 A discrete-time unit step function $u[k]$.Fig. 8.5 The λ -plane, the γ -plane and their mapping.

3. Discrete-Time Exponential γ^k

A continuous-time exponential $e^{\lambda t}$ can be expressed in an alternate form as

$$e^{\lambda t} = \gamma^t \quad (\gamma = e^\lambda \text{ or } \lambda = \ln \gamma) \quad (8.3a)$$

For example, $e^{-0.3t} = (0.7408)^t$ because $e^{-0.3} = 0.7408$. Conversely, $4^t = e^{1.386t}$ because $\ln 4 = 1.386$, that is, $e^{1.386} = 4$. In the study of continuous-time signals and systems we prefer the form $e^{\lambda t}$ rather than γ^t . The discrete-time exponential can also be expressed in two forms as

$$e^{\lambda k} = \gamma^k \quad (\gamma = e^\lambda \text{ or } \lambda = \ln \gamma) \quad (8.3b)$$

For example, $e^{3k} = (e^3)^k = (20.086)^k$. Similarly, $5^k = e^{1.609k}$ because $5 = e^{1.609}$. In the study of discrete-time signals and systems, unlike the continuous-time case, the form γ^k proves more convenient than the form $e^{\lambda k}$. Because of unfamiliarity with exponentials with bases other than e , exponentials of the form γ^k may seem inconvenient and confusing at first. The reader is urged to plot some exponentials to acquire a sense of these functions.

Nature of γ^k : The signal $e^{\lambda k}$ grows exponentially with k if $\text{Re } \lambda > 0$ (λ in RHP), and decays exponentially if $\text{Re } \lambda < 0$ (λ in LHP). It is constant or oscillates with constant amplitude if $\text{Re } \lambda = 0$ (λ on the imaginary axis). Clearly, the location of λ in the complex plane indicates whether the signal $e^{\lambda k}$ grows exponentially, decays exponentially, or oscillates with constant frequency (Fig. 8.5a). A constant signal ($\lambda = 0$) is also an oscillation with zero frequency. We now find a similar criterion for determining the nature of γ^k from the location of γ in the complex plane.

Figure 8.5a shows a complex plane (λ -plane). Consider a signal $e^{j\Omega k}$. In this case, $\lambda = j\Omega$ lies on the imaginary axis (Fig. 8.5a), and therefore is a constant-amplitude oscillating signal. This signal $e^{j\Omega k}$ can be expressed as γ^k , where $\gamma = e^{j\Omega}$.

disturbed, will neither go back to the original state nor continue to move farther away from the original state. The cone in this case is said to be in a **neutral equilibrium**.

Let us apply these observations to systems in general. If, in the absence of an external input, a system remains in a particular state (or condition) indefinitely, then that state is said to be an **equilibrium state of the system**. For an LTI system this equilibrium state is the zero state, in which all initial conditions are zero. Now suppose an LTI system is in equilibrium (zero state) and we change this state by creating some nonzero initial conditions. By analogy with the cone, if the system is stable it should eventually return to zero state. In other words, when left to itself, the system's output due to the nonzero initial conditions should approach 0 as $t \rightarrow \infty$. But the system output generated by initial conditions (zero-input response) is made up of its characteristic modes. For this reason we define stability as follows: a system is **(asymptotically) stable** if, and only if, all its characteristic modes $\rightarrow 0$ as $t \rightarrow \infty$. If any of the modes grows without bound as $t \rightarrow \infty$, the system is **unstable**. There is also a borderline situation in which the zero-input response remains bounded (approaches neither zero nor infinity), approaching a constant or oscillating with a constant amplitude as $t \rightarrow \infty$. For this borderline situation, the system is said to be **marginally stable** or just stable.

If an LTIC system has n distinct characteristic roots $\lambda_1, \lambda_2, \dots, \lambda_n$, the zero-input response is given by

$$y_0(t) = \sum_{j=1}^n c_j e^{\lambda_j t} \quad (2.62)$$

We have shown elsewhere [see Eq. (B.14)]

$$\lim_{t \rightarrow \infty} e^{\lambda t} = \begin{cases} 0 & \text{Re } \lambda < 0 \\ \infty & \text{Re } \lambda > 0 \end{cases} \quad (2.63)$$

It is helpful to study system stability in terms of the location of the system's characteristic roots in the complex plane. Let us first assume that the system has distinct roots only. If a characteristic root λ is located in the left half of the complex plane (LHP), its real part is negative ($\text{Re } \lambda < 0$). Similarly, if a root λ is located in the right half of the complex plane (RHP), its real part is positive ($\text{Re } \lambda > 0$). Along the imaginary axis, the real part is zero ($\text{Re } \lambda = 0$). These regions are delineated in Fig. 2.15. Equation (2.63) clearly shows that the characteristic modes corresponding to roots in LHP vanish as $t \rightarrow \infty$, while the modes corresponding to roots in RHP grow without bound as $t \rightarrow \infty$. However, the modes corresponding to simple (unrepeated) roots on the imaginary axis are of the form $e^{j\beta t}$; these are bounded (neither vanish nor grow without limit) as $t \rightarrow \infty$.

From this discussion it follows that a system is asymptotically stable if, and only if, all of its characteristic roots lie in the left half of the complex plane. If any of the roots—even one—lies in RHP, the system is unstable. If none of the roots lie in RHP, but if some unrepeated (simple) roots lie on the imaginary axis, then the system is marginally stable (Fig. 2.15).

So far we have assumed all of the system's n roots to be distinct. The modes corresponding to a root λ repeated r times are $e^{\lambda t}, t e^{\lambda t}, t^2 e^{\lambda t}, \dots, t^{r-1} e^{\lambda t}$. But as

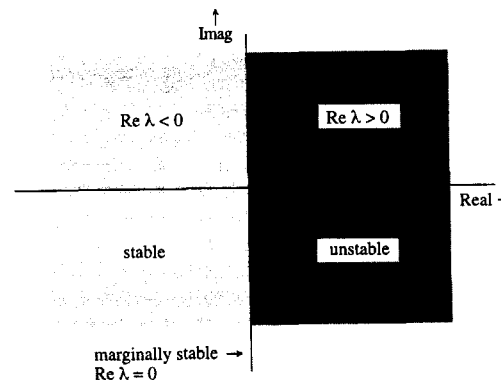


Fig. 2.15 Characteristic roots location and system stability.

$t \rightarrow \infty$, $t^k e^{\lambda t} \rightarrow 0$, if $\text{Re } \lambda < 0$ (λ in LHP). Therefore, repeated roots in LHP do not cause instability. But when the repeated roots are on the imaginary axis ($\lambda = j\omega$), the corresponding modes $t^k e^{j\omega t}$ approach infinity as $t \rightarrow \infty$. Therefore, repeated roots on the imaginary axis cause instability. Figure 2.16 shows characteristic modes corresponding to characteristic roots at various location in the complex plane. Observe the central role played by the characteristic roots or characteristic modes in determining the system's stability.

To summarize:

1. An LTIC system is asymptotically stable if, and only if, all the characteristic roots are in the LHP. The roots may be simple (unrepeated) or repeated.
2. An LTIC system is unstable if, and only if, either one or both of the following conditions exist: (i) at least one root is in the RHP, (ii) there are repeated roots on the imaginary axis.
3. An LTIC system is marginally stable if, and only if, there are no roots in the RHP, and there are some unrepeated roots on the imaginary axis.

Example 2.12

Investigate the stability of LTIC system described by the following equations:

- (a) $(D+1)(D^2+4D+8)y(t) = (D-3)f(t)$
- (b) $(D-1)(D^2+4D+8)y(t) = (D+2)f(t)$
- (c) $(D+2)(D^2+4)y(t) = (D^2+D+1)f(t)$
- (d) $(D+1)(D^2+4)^2y(t) = (D^2+2D+8)f(t)$

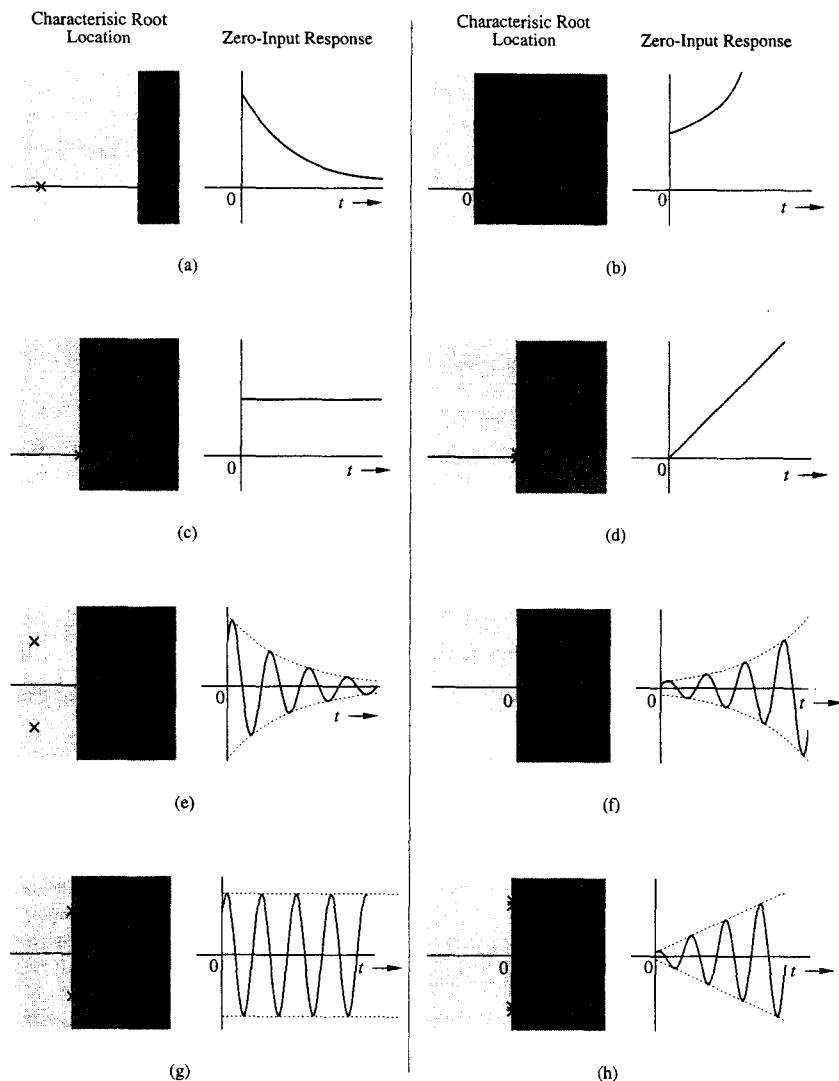


Fig. 2.16 Location of characteristic roots and the corresponding characteristic modes.

The characteristic polynomials of these systems are

$$(a) \quad (\lambda + 1)(\lambda^2 + 4\lambda + 8) = (\lambda + 1)(\lambda + 2 - j2)(\lambda + 2 + j2)$$

$$(b) \quad (\lambda - 1)(\lambda^2 + 4\lambda + 8) = (\lambda - 1)(\lambda + 2 - j2)(\lambda + 2 + j2)$$

$$(c) \quad (\lambda + 2)(\lambda^2 + 4) = (\lambda + 2)(\lambda - j2)(\lambda + j2)$$

$$(d) \quad (\lambda + 1)(\lambda^2 + 4)^2 = (\lambda + 2)(\lambda - j2)^2(\lambda + j2)^2$$

Consequently, the characteristic roots of the systems above are (see Fig. 2.17):

$$(a) \quad -1, -2 \pm j2 \quad (b) \quad 1, -2 \pm j2 \quad (c) \quad -2, \pm j2 \quad (d) \quad -1, \pm j2, \pm j2.$$

System (a) is asymptotically stable (all roots in LHP), (b) is unstable (one root in RHP), (c) is marginally stable (unrepeated roots on imaginary axis) and no roots in RHP, and (d) is unstable (repeated roots on the imaginary axis). ■

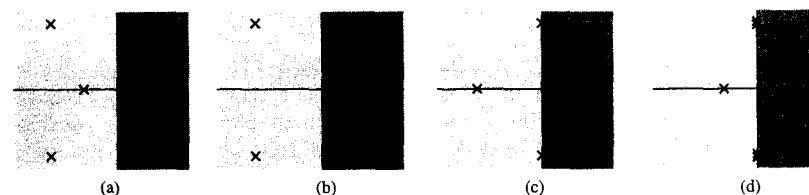


Fig. 2.17 Location of characteristic roots for systems in Example 2.12.

△ Exercise E2.16

For each of the systems specified by the equations below, plot its characteristic roots in the complex plane and determine whether it is asymptotically stable, marginally stable, or unstable.

$$(a) \quad D(D + 2)y(t) = 3f(t)$$

$$(b) \quad D^2(D + 3)y(t) = (D + 5)f(t)$$

$$(c) \quad (D + 1)(D + 2)y(t) = (2D + 3)f(t)$$

$$(d) \quad (D^2 + 1)(D^2 + 9)y(t) = (D^2 + 2D + 4)f(t)$$

$$(e) \quad (D + 1)(D^2 - 4D + 9)y(t) = (D + 7)f(t)$$

Answer: (a) marginally stable (b) unstable (c) stable (d) marginally stable (e) unstable. ▽

2.6-1 System Response to Bounded Inputs

From the example of the right circular cone, it appears that when a system is in stable equilibrium, application of a small force (input) produces a small response. In contrast, when the system is in unstable equilibrium, a small force (input) produces an unbounded response. Intuitively we feel that every bounded input should produce a bounded response in a stable system, whereas in an unstable system this would not be the case. We shall now verify this hunch and show that it is indeed true.

Recall that for an LTIC system

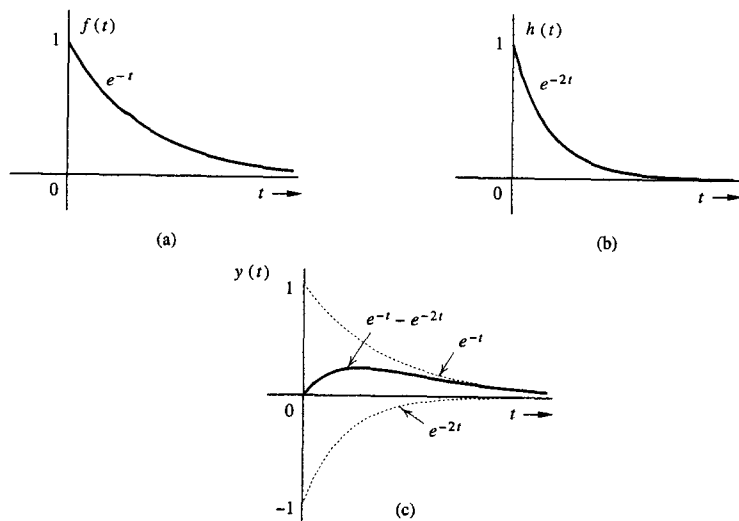
$$\begin{aligned} y(t) &= h(t) * f(t) \\ &= \int_{-\infty}^{\infty} h(\tau)f(t - \tau) d\tau \end{aligned} \quad (2.64)$$

Therefore

$$|y(t)| \leq \int_{-\infty}^{\infty} |h(\tau)||f(t - \tau)| d\tau$$

Moreover, if $f(t)$ is bounded, then $|f(t - \tau)| < K_1 < \infty$, and

$$|y(t)| \leq K_1 \int_{-\infty}^{\infty} |h(\tau)| d\tau$$

Fig. 2.6 Convolution of $f(t)$ and $h(t)$ in Example 2.4.

Because this integration is with respect to τ , we can pull e^{-2t} outside the integral, giving us

$$y(t) = e^{-2t} \int_0^t e^{\tau} d\tau = e^{-2t}(e^t - 1) = e^{-t} - e^{-2t} \quad t \geq 0 \quad (2.41)$$

Also, $y(t) = 0$ when $t < 0$ [see Eq. (2.38)]. This result, along with Eq. (2.41), yields

$$y(t) = (e^{-t} - e^{-2t})u(t) \quad (2.42)$$

The response is depicted in Fig. 2.6c. ■

△ Exercise E2.5

For an LTIC system with the impulse response $h(t) = 6e^{-t}u(t)$, determine the system response to the input: (a) $2u(t)$ and (b) $3e^{-3t}u(t)$.

Answer: (a) $12(1 - e^{-t})u(t)$ (b) $9(e^{-t} - e^{-3t})u(t)$ ▽

△ Exercise E2.6

Repeat Exercise E2.5 if the input $f(t) = e^{-t}u(t)$.

Answer: $6te^{-t}u(t)$ ▽

The Convolution Table

The task of convolution is considerably simplified by a ready-made convolution table (Table 2.1). This table, which lists several pairs of signals and their resulting convolution, can conveniently determine $y(t)$, a system response to an input $f(t)$, without performing the tedious job of integration. For instance, we could have readily found the convolution in Example 2.4 using pair 4 (with $\lambda_1 = -1$ and $\lambda_2 = -2$) to be $(e^{-t} - e^{-2t})u(t)$. The following example demonstrates the utility of this table.

TABLE 2.1: Convolution Table

No	$f_1(t)$	$f_2(t)$	$f_1(t) * f_2(t) = f_2(t) * f_1(t)$
1	$f(t)$	$\delta(t - T)$	$f(t - T)$
2	$e^{\lambda t}u(t)$	$u(t)$	$\frac{1 - e^{\lambda t}}{-\lambda} u(t)$
3	$u(t)$	$u(t)$	$tu(t)$
4	$e^{\lambda_1 t}u(t)$	$e^{\lambda_2 t}u(t)$	$\frac{e^{\lambda_1 t} - e^{\lambda_2 t}}{\lambda_1 - \lambda_2} u(t) \quad \lambda_1 \neq \lambda_2$
5	$e^{\lambda t}u(t)$	$e^{\lambda t}u(t)$	$te^{\lambda t}u(t)$
6	$te^{\lambda t}u(t)$	$e^{\lambda t}u(t)$	$\frac{1}{2}t^2 e^{\lambda t}u(t)$
7	$t^n u(t)$	$e^{\lambda t}u(t)$	$\frac{n! e^{\lambda t}}{\lambda^{n+1}} u(t) - \sum_{j=0}^n \frac{n! t^{n-j}}{\lambda^{j+1}(n-j)!} u(t)$
8	$t^m u(t)$	$t^n u(t)$	$\frac{m!n!}{(m+n+1)!} t^{m+n+1} u(t)$
9	$te^{\lambda_1 t}u(t)$	$e^{\lambda_2 t}u(t)$	$\frac{e^{\lambda_2 t} - e^{\lambda_1 t} + (\lambda_1 - \lambda_2)te^{\lambda_1 t}}{(\lambda_1 - \lambda_2)^2} u(t)$
10	$t^m e^{\lambda t}u(t)$	$t^n e^{\lambda t}u(t)$	$\frac{m!n!}{(n+m+1)!} t^{m+n+1} e^{\lambda t}u(t)$
11	$t^m e^{\lambda_1 t}u(t)$	$t^n e^{\lambda_2 t}u(t)$	$\sum_{j=0}^m \frac{(-1)^j m!(n+j)! t^{m-j} e^{\lambda_1 t}}{j!(m-j)!(\lambda_1 - \lambda_2)^{n+j+1}} u(t)$ $+ \sum_{k=0}^n \frac{(-1)^k n!(m+k)! t^{n-k} e^{\lambda_2 t}}{k!(n-k)!(\lambda_2 - \lambda_1)^{m+k+1}} u(t)$ $\lambda_1 \neq \lambda_2$
12	$e^{-\alpha t} \cos(\beta t + \theta)u(t)$	$e^{\lambda t}u(t)$	$\frac{\cos(\theta - \phi)e^{\lambda t} - e^{-\alpha t} \cos(\beta t + \theta - \phi)}{\sqrt{(\alpha + \lambda)^2 + \beta^2}} u(t)$ $\phi = \tan^{-1}[-\beta/(\alpha + \lambda)]$
13	$e^{\lambda_1 t}u(t)$	$e^{\lambda_2 t}u(-t)$	$\frac{e^{\lambda_1 t}u(t) + e^{\lambda_2 t}u(-t)}{\lambda_2 - \lambda_1} \quad \text{Re } \lambda_2 > \text{Re } \lambda_1$
14	$e^{\lambda_1 t}u(-t)$	$e^{\lambda_2 t}u(-t)$	$\frac{e^{\lambda_1 t} - e^{\lambda_2 t}}{\lambda_2 - \lambda_1} u(-t)$

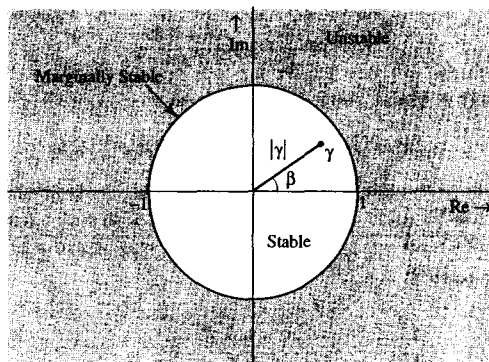


Fig. 9.6 Characteristic roots location and system stability.

It is clear that a system is asymptotically stable if and only if

$$|\gamma_i| < 1 \quad i = 1, 2, \dots, n$$

These results can be grasped more effectively in terms of the location of characteristic roots in the complex plane. Figure 9.6 shows a circle of unit radius, centered at the origin in a complex plane. Our discussion clearly shows that if all characteristic roots of the system lie inside this circle (**unit circle**), $|\gamma_i| < 1$ for all i and the system is asymptotically stable. On the other hand, even if one characteristic root lies outside the unit circle, the system is unstable. If none of the characteristic roots lie outside the unit circle, but some simple (unrepeated) roots lie on the circle itself, the system is marginally stable. If two or more characteristic roots coincide on the unit circle (repeated roots), the system is unstable. The reason is that for repeated roots, the zero-input response is of the form $k^{r-1}\gamma^k$, and if $|\gamma| = 1$, then $|k^{r-1}\gamma^k| = k^{r-1} \rightarrow \infty$ as $k \rightarrow \infty$.† Note, however, that repeated roots inside the unit circle do not cause instability. Figure 9.7 shows the characteristic modes corresponding to characteristic roots at various locations in the complex plane. To summarize:

1. An LTID system is asymptotically stable if and only if all the characteristic roots are inside the unit circle. The roots may be simple or repeated.
2. An LTID system is unstable if and only if either one or both of the following conditions exist: (i) at least one root is outside the unit circle; (ii) there are repeated roots on the unit circle.
3. An LTID system is marginally stable if and only if there are no roots outside the unit circle and there are some unrepeated roots on the unit circle.

†If the development of discrete-time systems is parallel to that of continuous-time systems, we wonder why the parallel breaks down here. Why, for instance, aren't LHP and RHP the regions demarcating stability and instability? The reason lies in the form of the characteristic modes. In continuous-time systems we chose the form of characteristic mode as $e^{\lambda_i t}$. In discrete-time systems we choose the form (for computational convenience) to be γ_i^k . Had we chosen this form to be $e^{\lambda_i k}$ where $\gamma_i = e^{\lambda_i}$, then LHP and RHP (for the location of λ_i) again would demarcate stability and instability. The reason is that if $\gamma = e^{\lambda}$, $|\gamma| = 1$ implies $|e^{\lambda}| = 1$, and therefore $\lambda = j\omega$. This shows that the unit circle in γ plane maps into the imaginary axis in the λ plane.

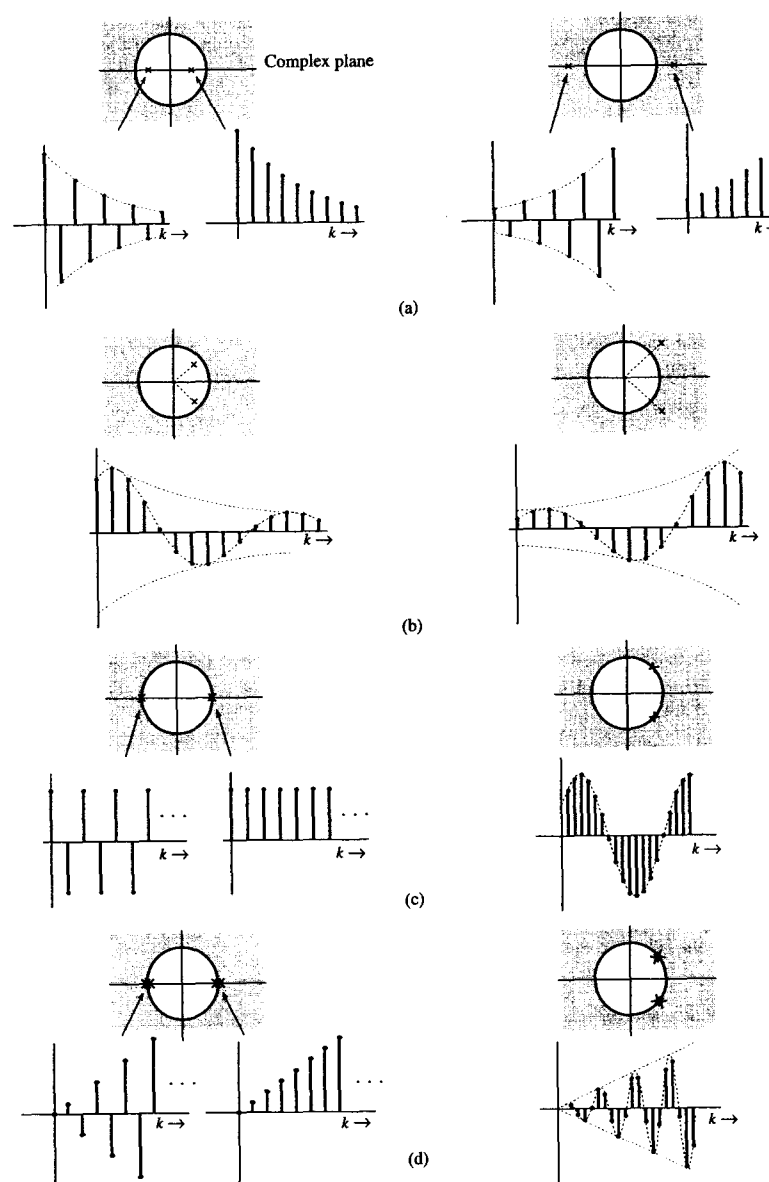


Fig. 9.7 Characteristic roots location and the corresponding characteristic modes.

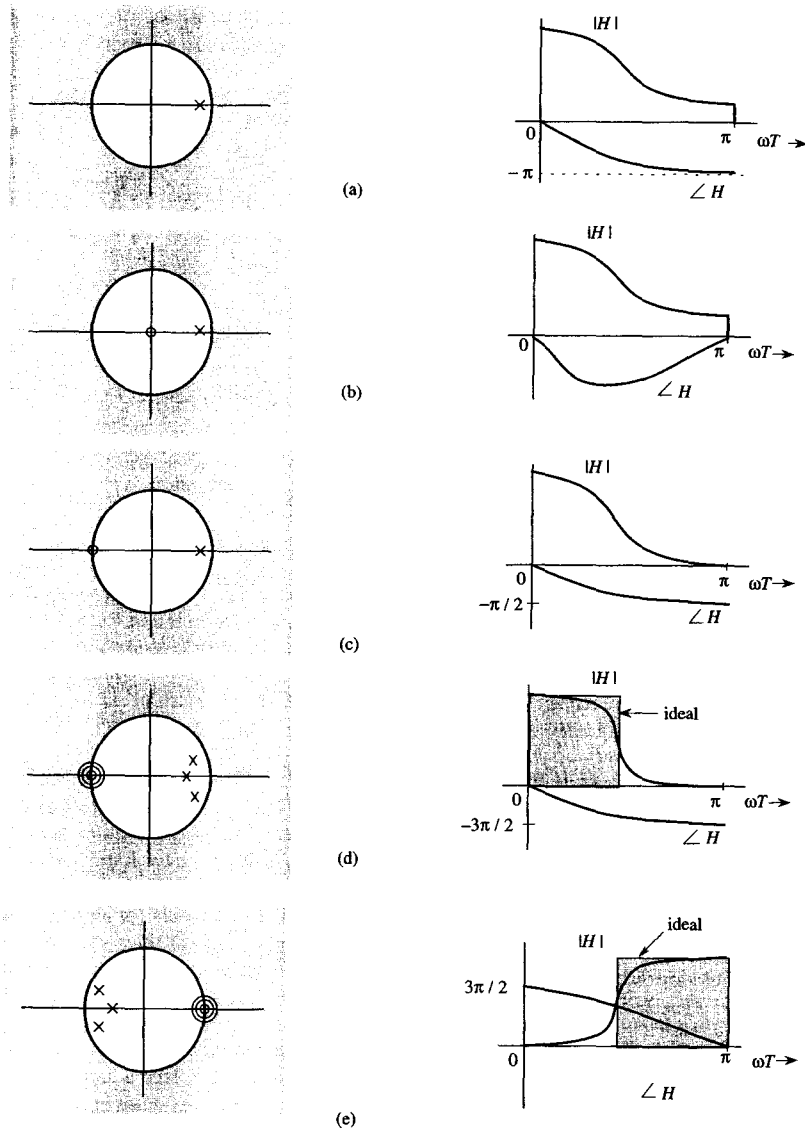


Fig. 12.4 Various pole-zero configurations and the corresponding frequency response.

$\angle H[e^{j\omega T}]$. The phase spectrum $-\omega T$ is a linear function of frequency and therefore represents a pure time-delay of T seconds (see Eq. (10.48) or Exercise E12.2). Therefore, a pole (a zero) at the origin causes a time delay (time advance) of T seconds in the response. There is no change in the amplitude response.

For a stable system, all the poles must be located inside the unit circle. The zeros may lie anywhere. Also, for a physically realizable system, $H[z]$ must be a proper fraction, that is, $n \geq m$. If, to achieve a certain amplitude response, we require $m > n$, we can still make the system realizable by placing a sufficient number of poles at the origin. This will not change the amplitude response but it will increase the time delay of the response.

In general, a pole at a point has the opposite effect of a zero at that point. Placing a zero closer to a pole tends to cancel the effect of that pole on the frequency response.

Lowpass Filters

A lowpass filter has a maximum gain at $\omega = 0$, which corresponds to point $e^{j0T} = 1$ on the unit circle. Clearly, placing a pole inside the unit circle near the point $z = 1$ (Fig. 12.4a) would result in a lowpass response. The corresponding amplitude and phase response appears in Fig. 12.4a. For smaller values of ω , the point $e^{j\omega T}$ (a point on the unit circle at an angle ωT) is closer to the pole, and consequently the gain is higher. As ω increases, the distance of the point $e^{j\omega T}$ from the pole increases. Consequently the gain decreases, resulting in a lowpass characteristic. Placing a zero at the origin does not change the amplitude response but it does modify the phase response, as illustrated in Fig. 12.4b. Placing a zero at $z = -1$, however, changes both the amplitude and phase response (Fig. 12.4c). The point $z = -1$ corresponds to frequency $\omega = \pi/T$ ($z = e^{j\omega T} = e^{j\pi} = -1$). Consequently, the amplitude response now becomes more attenuated at higher frequencies, with a zero gain at $\omega T = \pi$. We can approach ideal lowpass characteristics by using more poles staggered near $z = 1$ (but within the unit circle). Figure 12.4d shows a third-order lowpass filter with three poles near $z = 1$ and a third-order zero at $z = -1$, with corresponding amplitude and phase response. For an ideal lowpass filter we need an enhanced gain at every frequency in the band $(0, \omega_c)$. This can be achieved by placing a continuous wall of poles (requiring an infinite number of poles) opposite this band.

Highpass Filters

A highpass filter has a small gain at lower frequencies and a high gain at higher frequencies. Such a characteristic can be realized by placing a pole or poles near $z = -1$ because we want the gain at $\omega T = \pi$ to be the highest. Placing a zero at $z = 1$ further enhances suppression of gain at lower frequencies. Figure 12.4e shows a possible pole-zero configuration of the third-order highpass filter with corresponding amplitude and phase response.

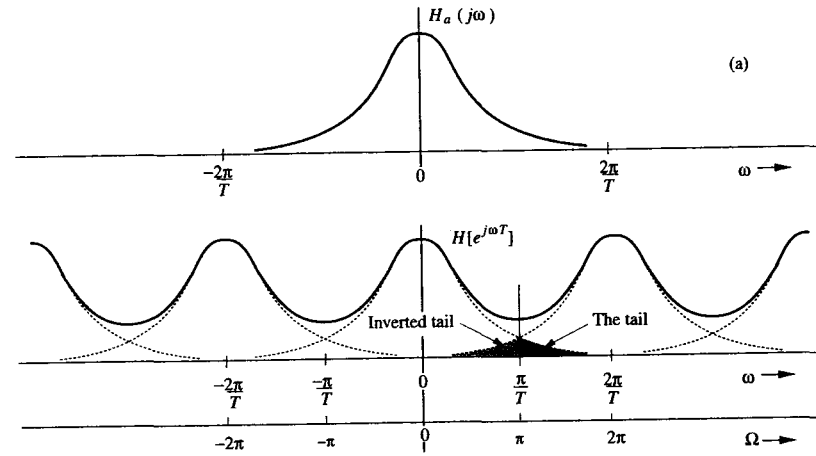
Example 12.2: Bandpass Filter

Using trial-and-error, design a tuned (bandpass) filter with zero transmission at 0 Hz and also at 500 Hz. The resonant frequency is required to be 125 Hz. The highest frequency to be processed is $F_h = 500$ Hz.

Because $F_h = 500$, we require $T \leq \frac{1}{1000}$ [see Eq. (8.17)]. Let us select $T = 10^{-3}$. Since the amplitude response is zero at $\omega = 0$ and $\omega = 1000\pi$, we need to place zeros at $e^{j\omega T}$ corresponding to $\omega = 0$ and $\omega = 1000\pi$. For $\omega = 0$, $z = e^{j\omega T} = 1$; for $\omega = 1000\pi$ (with $T = 10^{-3}$), $e^{j\omega T} = -1$. Hence, there must be zeros at $z = \pm 1$. Moreover, we need enhanced frequency response at $\omega = 250\pi$. This frequency (with $\omega T = \pi/4$) corresponds to $z = e^{j\omega T} = e^{j\pi/4}$. Therefore, to enhance the frequency response at this frequency, we

TABLE 12.1

	$H_a(s)$	$h_a(t)$	$h[k]$	$H[z]$
1	K	$K\delta(t)$	$TK\delta[k]$	TK
2	$\frac{1}{s}$	$u(t)$	$Tu[k]$	$\frac{Tz}{z-1}$
3	$\frac{1}{s^2}$	t	kT^2	$\frac{T^2z}{(z-1)^2}$
4	$\frac{1}{s^3}$	$\frac{t^2}{2}$	$\frac{k^2T^3}{2}$	$\frac{T^3z(z+1)}{2(z-1)^3}$
5	$\frac{1}{s-\lambda}$	$e^{\lambda t}$	$Te^{\lambda kT}$	$\frac{Tz}{z-e^{\lambda T}}$
6	$\frac{1}{(s-\lambda)^2}$	$te^{\lambda t}$	$kT^2e^{\lambda kT}$	$\frac{T^2ze^{\lambda T}}{(z-e^{\lambda T})^2}$
7	$\frac{As+B}{s^2+2as+c}$ $r = \sqrt{\frac{A^2c+B^2-2ABa}{c-a^2}}$	$Tr e^{-at} \cos(bt+\theta)$ $b = \sqrt{c-a^2}$ $\theta = \tan^{-1} \left(\frac{Aa-B}{A\sqrt{c-a^2}} \right)$	$Tr e^{-\alpha kT} \cos(bkT+\theta)$	$\frac{Trz [z \cos \theta - e^{-aT} \cos(bT-\theta)]}{z^2 - (2e^{-aT} \cos bT)z + e^{-2aT}}$

Fig. 12.10 Aliasing in digital filters, and a choice of the sampling interval T .

$$h[k] = \lim_{T \rightarrow 0} Th_a(kT)$$

In Chapter 5 (Fig. 5.6), we showed that the Fourier transform of the samples of $h_a(t)$ consists of periodic repetition of $H_a(j\omega)$ with period equal to the sampling frequency $\omega_s = 2\pi/T = 2\pi\mathcal{F}_s$.[†] Also $H_a(j\omega)$ is not generally bandlimited. Hence, aliasing among various repeating cycles cannot be prevented, as depicted in Fig. 12.10b. The resulting spectrum will be different from the desired spectrum, especially at higher frequencies. If $H_a(j\omega)$ were to be bandlimited; that is, if $H_a(j\omega) = 0$ for $|\omega| > \omega_0$, then the overlap could be avoided if we select the period $2\pi/T > 2\omega_0$. However, according to the Paley-Wiener criterion [Eq. (4.61)], every practical system frequency response is nonbandlimited, and the cycle overlap is inevitable. However, for frequencies beyond some ω_0 , if $|H_a(j\omega)|$ is a negligible fraction, say 1%, of $|H_a(j\omega)|_{\max}$, then we can consider¹ $H_a(j\omega)$ to be essentially bandlimited to ω_0 , and we can select

$$T = \frac{\pi}{\omega_0} \quad (12.46)$$

Example 12.4

Design a digital filter to realize the first-order analog lowpass Butterworth filter with the transfer function

$$H_a(s) = \frac{\omega_c}{s + \omega_c} \quad \omega_c = 10^5 \quad (12.47)$$

[†]How can we apply the discussion in Chapter 5, which applies to impulse samples of continuous-time signals, to discrete-time signals? Recall our discussion in Sec. 10.4 (Fig. 10.8), where we showed that the spectrum of discrete-time signal is just a scaled version of the spectrum of the impulse samples of the corresponding continuous-time signal.

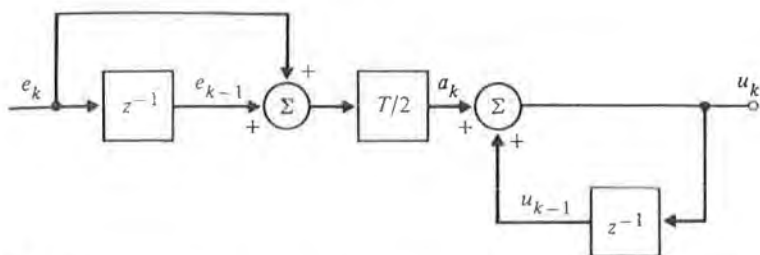


Figure 2.4 A block diagram of trapezoid integration as represented by (2.7).

value of the integral estimate, u_k . The discrete integration occurs in the loop with one delay, z^{-1} , and unity gain.

2.3.3 Block Diagrams and State-Variable Descriptions

Because (2.16) is a linear algebraic relationship, a system of such relations is described by a system of linear equations. These can be solved by the methods of linear algebra or by the graphical methods of block diagrams. To use block-diagram analysis to manipulate these discrete-transfer-function relationships, there are only four primitive cases:

1. The transfer function of paths in parallel is the sum of the single-path transfer functions (Fig. 2.5).
2. The transfer function of paths in series is the *product* of the path transfer functions (Fig. 2.6).
3. The transfer function of a single loop of paths is the transfer function of the forward path divided by one minus the loop transfer function (Fig. 2.7).
4. The transfer function of an arbitrary multipath diagram is given by combinations of these cases. Mason's rule⁶ can also be used.

For the general difference equation of (2.2), we already have the transfer function in (2.15). It is interesting to connect this case with a block diagram *using only simple delay forms for z* in order to see several "canonical" block diagrams and to introduce the description of discrete systems using equations of state.

There are many ways to reduce the difference equation (2.2) to a block diagram involving z only as the delay operator, z^{-1} . The first one we will

⁶Mason (1956). See Franklin, Powell, and Emami-Naeini (1986) for a discussion.

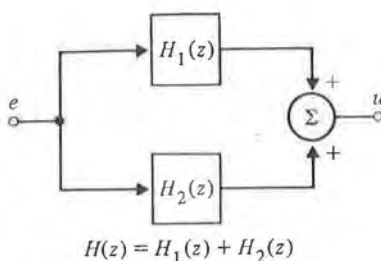


Figure 2.5 Block diagram of parallel blocks.

consider leads to the “control” canonical form. We begin with the transfer function as a ratio of polynomials

$$U(z) = H(z)E(z) = \frac{b(z)}{a(z)}E(z) = b(z)\xi,$$

where

$$\xi = \frac{E(z)}{a(z)}$$

and thus

$$a(z)\xi = E(z).$$

At this point we need to get specific; and rather than carry through with a system of arbitrary order, we will work out the details for the third-order case and leave it to the reader to extend the results in the obvious way to whatever order is desired. In the development that follows, we will consider the variables u , e , and ξ as *time* variables and z as an advance operator such that $zu(k) = u(k+1)$ or $z^{-1}u = u(k-1)$. With this convention (which is simply using the property of z derived earlier), consider the equations

$$(z^3 + a_1z^2 + a_2z + a_3)\xi = e, \quad (2.18)$$

$$(b_0z^3 + b_1z^2 + b_2z + b_3)\xi = u. \quad (2.19)$$

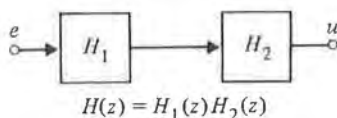


Figure 2.6 Block diagram of cascade blocks.

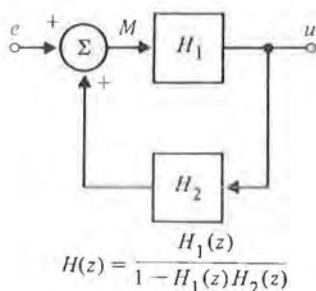


Figure 2.7 Feedback transfer function.

We can write (2.18) as

$$z^3\xi = e - a_1z^2\xi - a_2z\xi - a_3\xi$$

$$\xi(k+3) = e(k) - a_1\xi(k+2) - a_2\xi(k+1) - a_3\xi(k). \quad (2.20)$$

Now assume we have $z^3\xi$, which is to say that we have $\xi(k+3)$ because z^3 is an advance operator of three steps. If we operate on this with z^{-1} three times in a row, we will get back to $\xi(k)$, as shown in Fig. 2.8(a). From (2.20), we can now compute $z^3\xi$ from e and the lower powers of z and ξ given in the block diagram; the picture is now as given in Fig. 2.8(b). To complete the representation of (2.18) and (2.19), we need only add the formation of the output u as a weighted sum of the variables $z^3\xi$, $z^2\xi$, $z\xi$, and ξ according to (2.19). The completed picture is shown in Fig. 2.8(c).

In Fig 2.8(c), the internal variables have been named x_1 , x_2 , and x_3 . These variables comprise the *state* of this dynamic system in this form. Having the block diagram shown in Fig. 2.8(c), we can write down, almost by inspection, the difference equations that describe the evolution of the state, again using the fact that the transfer function z^{-1} corresponds to a one-unit delay. For example, we see that $x_3(k+1) = x_2(k)$ and $x_2(k+1) = x_1(k)$. Finally, expressing the sum at the far left of the figure, we have

$$x_1(k+1) = -a_1x_1(k) - a_2x_2(k) - a_3x_3(k) + e(k).$$

We collect these three equations together in proper order, and we have

$$\begin{aligned} x_1(k+1) &= -a_1x_1(k) - a_2x_2(k) - a_3x_3(k) + e(k), \\ x_2(k+1) &= x_1(k), \\ x_3(k+1) &= x_2(k). \end{aligned} \quad (2.21)$$

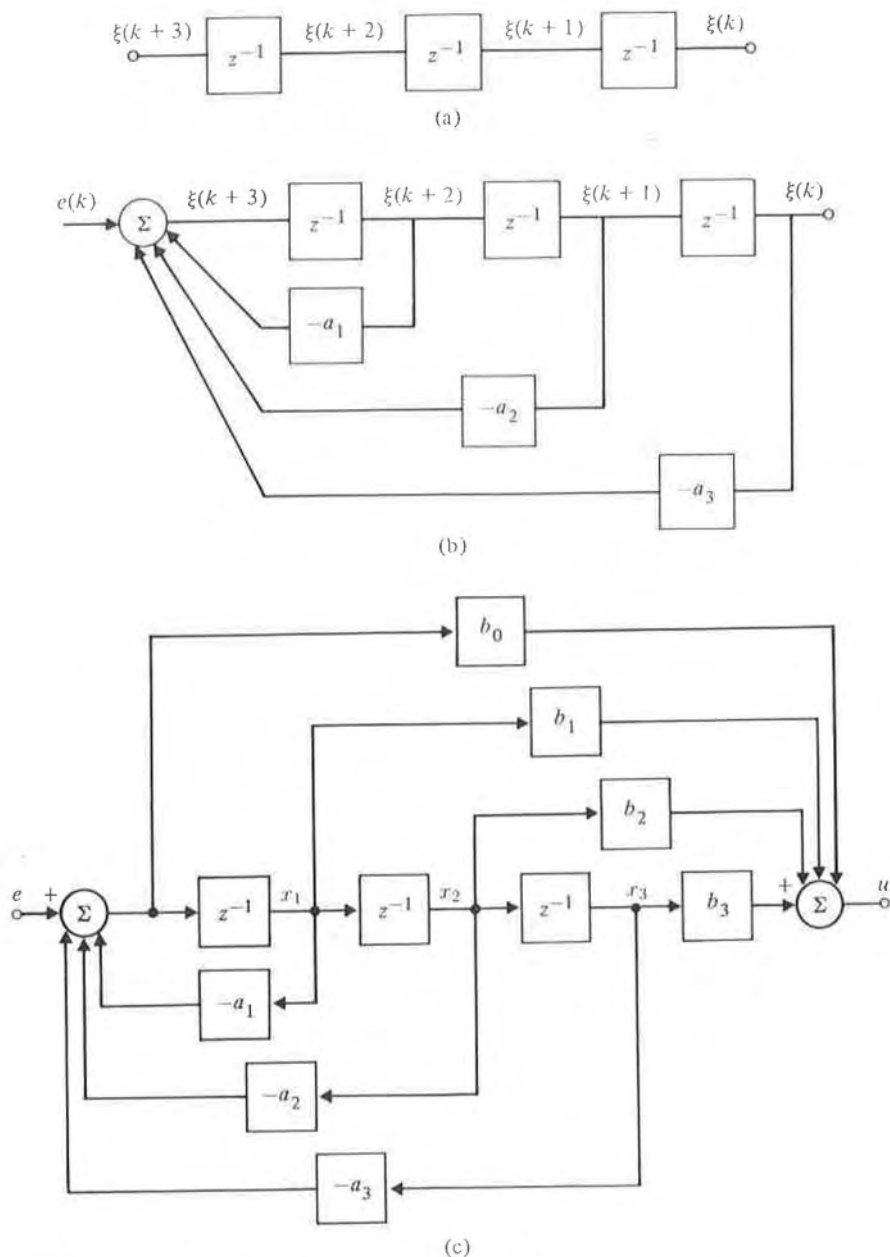


Figure 2.8 Block diagram development of control canonical form. (a) Solving for $\xi(k)$; (b) solving for $\xi(k+3)$ from $e(k)$ and past ξ 's; (c) solving for $U(k)$ from ξ 's.

Using vector-matrix notation,⁷ we can write this in the compact form

$$\mathbf{x}(k+1) = \mathbf{A}_c \mathbf{x}(k) + \mathbf{B}_c e(k),$$

where

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix},$$

$$\mathbf{A}_c = \begin{bmatrix} -a_1 & -a_2 & -a_3 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \quad (2.22a)$$

and

$$\mathbf{B}_c = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad (2.22b)$$

The output equation is also immediate except that we must watch to catch *all* paths by which the state variables combine in the output. The problem is caused by the b_0 term. If $b_0 = 0$, then $u = b_1 x_1 + b_2 x_2 + b_3 x_3$, and the corresponding matrix form is immediate. However, if b_0 is not 0, x_1 for example not only reaches the output through b_1 but also by the parallel path with gain $-b_0 a_1$. The complete equation is

$$u = (b_1 - a_1 b_0)x_1 + (b_2 - a_2 b_0)x_2 + (b_3 - a_3 b_0)x_3 + b_0 e.$$

In vector/matrix notation, we have

$$u = \mathbf{C}_c \mathbf{x} + \mathbf{D}_c e,$$

where

$$\mathbf{C}_c = [b_1 - a_1 b_0 \quad b_2 - a_2 b_0 \quad b_3 - a_3 b_0], \quad (2.23a)$$

$$\mathbf{D}_c = [b_0]. \quad (2.23b)$$

⁷We assume the reader has some knowledge of matrices. The results we require and references to study material are given in Appendix C. To distinguish vectors and matrices from scalar variables, we will use bold-face type.

2.5.1 The Unit Pulse

We have already seen that the unit pulse is defined by²¹

$$\begin{aligned} e_1(k) &= 1 & (k = 0) \\ &= 0 & (k \neq 0) \\ &= \delta_k; \end{aligned}$$

therefore we have

$$E_1(z) = \sum_{k=-\infty}^{\infty} \delta_k z^{-k} = z^0 = 1. \quad (2.80)$$

This result is much like the continuous case, wherein the Laplace transform of the unit impulse is the constant 1.0.

The quantity $E_1(z)$ gives us an instantaneous method to relate signals to systems: To characterize the system $H(z)$, consider the signal $u(k)$, which is the unit pulse response; then $U(z) = H(z)$.

2.5.2 The Unit Step

Consider the unit step function defined by

$$\begin{aligned} e_2(k) &= 1 & (k \geq 0) \\ &= 0 & (k < 0) \\ &\triangleq 1(k). \end{aligned}$$

In this case, the z -transform is

$$\begin{aligned} E_2(z) &= \sum_{k=-\infty}^{\infty} e_2(k) z^{-k} = \sum_{k=0}^{\infty} z^{-k} \\ &= \frac{1}{1 - z^{-1}} & (|z^{-1}| < 1) \\ &= \frac{z}{z - 1} & (|z| > 1). \end{aligned} \quad (2.81)$$

²¹We have shifted notation here to use $e(k)$ rather than e_k for the k th sample. We use subscripts to identify different signals.

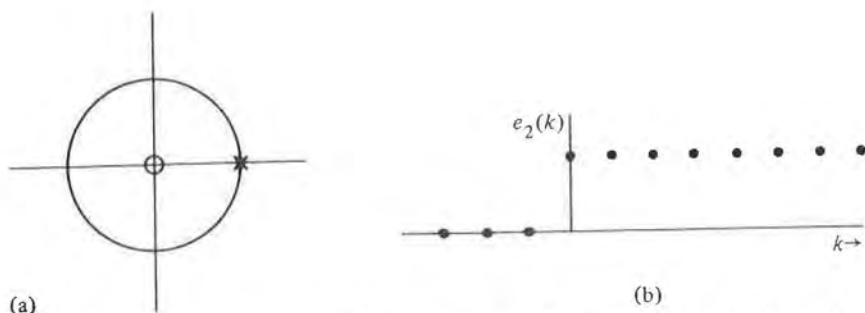


Figure 2.22 (a) Pole and zero of $E_2(z)$ in the z -plane. The unit circle is shown for reference. (b) Plot of $e_2(k)$.

Here the transform is characterized by a zero at $z = 0$ and a pole at $z = 1$. The significance of the convergence being restricted to $|z| > 1$ will be explored later when we consider the inverse transform operation. The Laplace transform of the unit step is $1/s$; we may thus keep in mind that a pole at $s = 0$ for a continuous signal corresponds in some way to a pole at $z = 1$ for discrete signals. We will explore this further later. In any event, we record that a pole at $z = 1$ with convergence outside the unit circle, $|z| = 1$, will correspond to a constant for positive time and zero for negative time.

To emphasize the connection between the time domain and the z -plane, we sketch in Fig. 2.22 the z -plane with the unit circle shown and the pole of $E_2(z)$ marked \times and the zero marked \circ . Beside the z -plane, we sketch the time plot of $e_2(k)$.

2.5.3 Exponential

The one-sided exponential in time is

$$\begin{aligned} e_3(k) &= r^k & (k \geq 0) \\ &= 0 & (k < 0), \end{aligned} \quad (2.82)$$

which is the same as $r^k 1(k)$, using the symbol $1(k)$ for the unit step function.

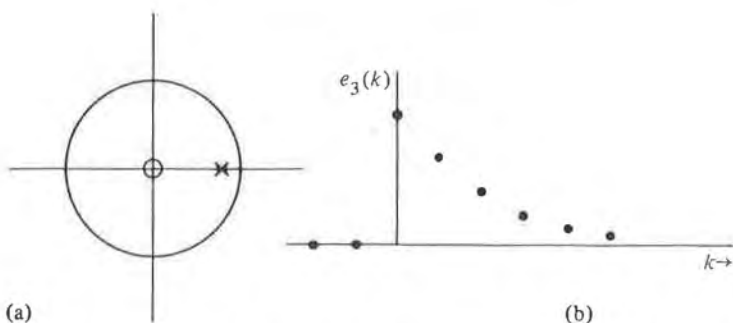


Figure 2.23 (a) Pole and zero of $E_3(z)$ in the z -plane. (b) Plot of $e_3(k)$.

Now we get

$$\begin{aligned}
 E_3(z) &= \sum_{k=0}^{\infty} r^k z^{-k} \\
 &= \sum_{k=0}^{\infty} (rz^{-1})^k \\
 &= \frac{1}{1 - rz^{-1}} \quad (|rz^{-1}| < 1) \\
 &= \frac{z}{z - r} \quad (|z| > |r|). \quad (2.83)
 \end{aligned}$$

The pole of $E_3(z)$ is at $z = r$. From (2.82) we know that $e_3(k)$ grows without bound if $|r| > 1$. From (2.83) we conclude that a z -transform that converges for large z and has a real pole *outside* the circle $|z| = 1$ corresponds to a growing signal. If such a signal were the unit-pulse response of our system, such as our digital control program, we would say the program was *unstable* as we saw in (2.37). We plot in Fig. 2.23 the z -plane and the corresponding time history of $E_3(z)$ as $e_3(k)$ for the stable value, $r = 0.6$.

2.5.4 General Sinusoid

Our next example considers the modulated sinusoid $e_4(k) = [r^k \cos k\theta]1(k)$, where we assume $r > 0$. Actually, we can decompose $e_4(k)$ into the sum of two complex exponentials as

$$e_4(k) = r^k \left(\frac{e^{jk\theta} + e^{-jk\theta}}{2} \right) 1(k),$$

and because the z -transform is linear,²² we need only compute the transform of each single complex exponential and add the results later. We thus take first

$$e_5(k) = r^k e^{jk\theta} 1(k) \quad (2.84)$$

and compute

$$\begin{aligned} E_5(z) &= \sum_{k=0}^{\infty} r^k e^{jk\theta} z^{-k} \\ &= \sum_{k=0}^{\infty} (r e^{j\theta} z^{-1})^k \\ &= \frac{1}{1 - r e^{j\theta} z^{-1}} \\ &= \frac{z}{z - r e^{j\theta}} \quad (|z| > r). \end{aligned} \quad (2.85)$$

The signal $e_5(k)$ grows without bound as k gets large if and only if $r > 1$, and a system with this pulse response is BIBO stable if and only if $|r| < 1$. The boundary of stability is the unit circle. To complete the argument given above for $e_4(k) = r^k \cos k\theta 1(k)$, we see immediately that the other half is found by replacing θ by $-\theta$ in (2.85),

$$\mathcal{Z}\{r^k e^{-jk\theta} 1(k)\} = \frac{z}{z - r e^{-j\theta}} \quad (|z| > r), \quad (2.86)$$

and thus that

$$\begin{aligned} E_4(z) &= \frac{1}{2} \left\{ \frac{z}{z - r e^{j\theta}} + \frac{z}{z - r e^{-j\theta}} \right\} \\ &= \frac{z(z - r \cos \theta)}{z^2 - 2r(\cos \theta)z + r^2} \quad (|z| > r). \end{aligned} \quad (2.87)$$

The z -plane pole-zero pattern of $E_4(z)$ and the time plot of $e_4(k)$ are shown in Fig. 2.24 for $r = 0.7$ and $\theta = 45^\circ$.

We note in passing that if $\theta = 0$, then e_4 reduces to e_3 and, with $r = 1$, to e_2 , so that three of our signals are special cases of e_4 . By exploiting the

²²We have not shown this formally. The demonstration, using the definition of linearity given above, is simple and is given in Section 2.7.

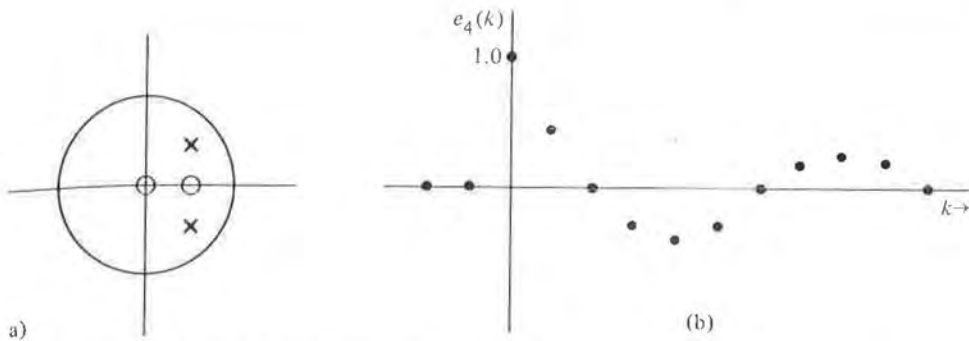


Figure 2.24 (a) Poles and zeros of $E_4(z)$ for $\theta = 45^\circ$, $r = 0.7$ in the z -plane. (b) Plot of $e_4(k)$.

atures of $E_4(z)$, we can draw a number of conclusions about the relation between pole locations in the z -plane and the time-domain signals to which the poles correspond. We collect these for later reference.

The settling time of a transient, defined as the time required for the signal to decay to one percent of its maximum value, is set mainly by the value of the radius, r , of the poles.

- $r > 1$ corresponds to a growing signal that will not decay at all.
- $r = 1$ corresponds to a signal with constant amplitude (which is *not* BIBO stable as a pulse response).
- For $r < 1$, the closer r is to 0 the shorter the settling time. The corresponding system is BIBO stable. We can compute the settling time in samples, N , in terms of the pole radius, r .

pole radius, r	response duration, N
0.9	43
0.8	21
0.6	9
0.4	5

- A pole at $r = 0$ corresponds to a transient of finite duration.

The number of samples per oscillation of a sinusoidal signal is determined by θ . If we require $\cos \theta k = \cos(\theta(k + N))$, we find that a period

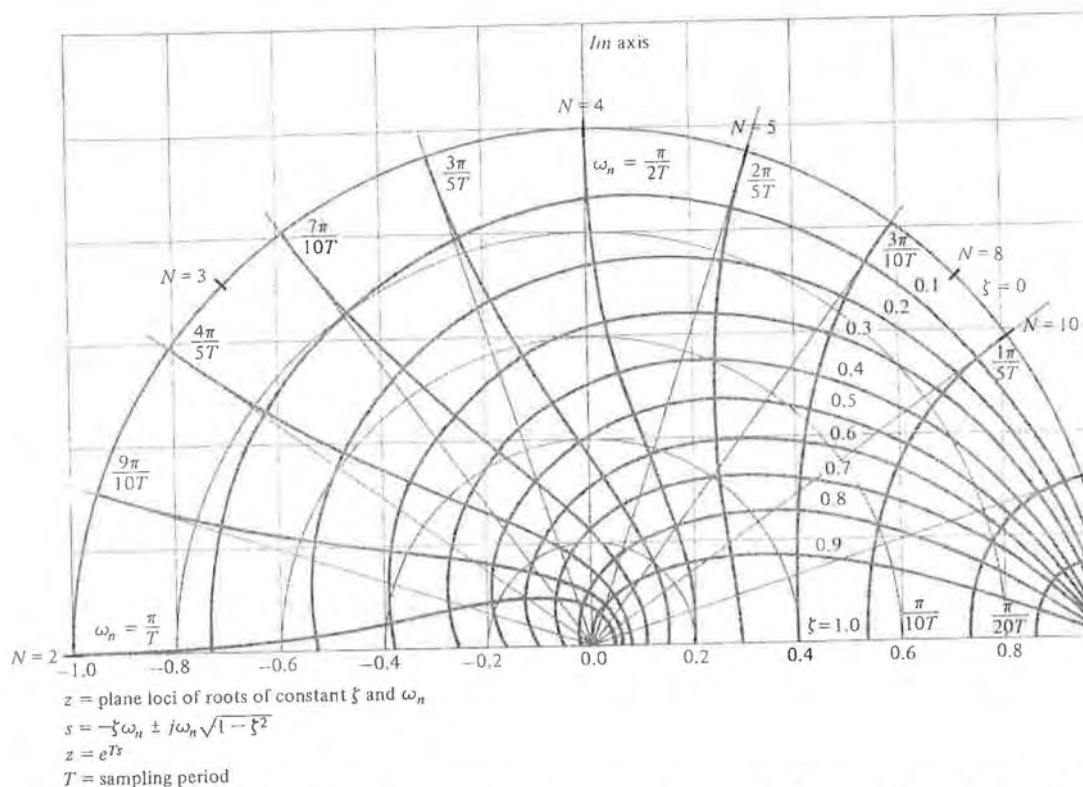


Figure 2.25 Sketch of the unit circle with angle θ marked in numbers of samples per cycle.

of 2π rad contains N samples, where

$$N = \left. \frac{2\pi}{\theta} \right|_{\text{rad}} = \left. \frac{360}{\theta} \right|_{\text{deg}} \text{ samples/cycle.}$$

For $\theta = 45^\circ$, we have $N = 8$, and the plot of $e_4(k)$ given in Fig. 2.24(b) shows the eight samples in the first cycle very clearly. A sketch of the unit circle with several points corresponding to various numbers of samples per cycle marked is drawn in Fig. 2.25. The sampling frequency in Hertz is $1/T$, and the signal frequency is $f = 1/NT$ so that $N = f_s/f$ and $1/N$ is a *normalized* signal frequency. Since $\theta = (2\pi)/N$, θ is the normalized signal frequency in radians/sample. θ/T is the frequency in radians/second.

2.5.5 Correspondence with Continuous Signals

From the calculation of these few z -transforms, we have established that the duration of a time signal is related to the radius of the pole locations and the number of samples per cycle is related to the angle, θ . Another set of very useful relationships can be established by considering the signals to be samples from a continuous signal, $e(t)$, with Laplace transform $E(s)$. With this device we can exploit our knowledge of s -plane features by transferring them to equivalent z -plane properties. For the specific numbers represented in the illustration of e_4 , we take the continuous signal

$$y(t) = e^{-at} \cos bt \, 1(t) \quad (2.88)$$

with

$$aT = 0.3567,$$

$$bT = \pi/4.$$

And, taking samples one second apart ($T = 1$), we have

$$\begin{aligned} y(kT) &= (e^{-0.3567})^k \cos \frac{\pi k}{4} 1(k) \\ &= (0.7)^k \cos \frac{\pi k}{4} 1(k) \\ &= e_4(k). \end{aligned}$$

The poles of the Laplace transform of $y(t)$ (in the s -plane) are at

$$s_{1,2} = -a + jb, -a - jb.$$

From (2.87), the z -transform of $E_4(z)$ has poles at

$$z_{1,2} = re^{j\theta}, re^{-j\theta},$$

but because $y(kT)$ equals $e_4(k)$, it follows that

$$\begin{aligned} r &= e^{-aT}, & \theta &= bT, \\ z_{1,2} &= e^{s_1 T}, & e^{s_2 T}. \end{aligned}$$

If $E(z)$ is a ratio of polynomials in z , which will be the case if $e(k)$ is generated by a linear difference equation with constant coefficients, then by partial fraction expansion, $E(z)$ can be expressed as a sum of elementary

B.1 PROPERTIES OF z-TRANSFORMS

Let $\mathcal{F}_i(s)$ be the Laplace transform of $f_i(t)$ and $F_i(z)$ be the z -transform of $f_i(kT)$.

Table B.1

Number	Laplace Transform	Samples	z -Transform	Comment
—	$\mathcal{F}_i(s)$	$f_i(kT)$	$F_i(z)$	
1	$\alpha\mathcal{F}_1(s) + \beta\mathcal{F}_2(s)$	$\alpha f_1(kT) + \beta f_2(kT)$	$\alpha F_1(z) + \beta F_2(z)$	The z -transform is linear
2	$\mathcal{F}_1(e^{Ts})\mathcal{F}_2(s)$	$\sum_{\ell=-\infty}^{\infty} f_1(\ell T)f_2(kT - \ell T)$	$F_1(z)F_2(z)$	Discrete convolution corresponds to product of z -transforms
3	$e^{+nTs}\mathcal{F}(s)$	$f(kT + nT)$	$z^n F(z)$	Shift in time
4	$\mathcal{F}(s + a)$	$e^{-akT}f(kT)$	$F(e^{aT}z)$	Shift in frequency
5	—	$\lim_{k \rightarrow \infty} f(kT)$	$\lim_{z \rightarrow 1} (z - 1)F(z)$	If all poles of $(z - 1)F(z)$ are inside the unit circle and $F(z)$ converges for $1 \leq z $
6	$\mathcal{F}(s/\omega_n)$	$f(\omega_n kT)$	$F(z; \omega_n T)$	Time and frequency scaling
7	—	$f_1(kT)f_2(kT)$	$\frac{1}{2\pi j} \oint_{c_3} F_1(\zeta)F_2(z/\zeta) \frac{d\zeta}{\zeta}$	Time product
8	$\mathcal{F}_3(s) = \mathcal{F}_1(s)\mathcal{F}_2(s)$	$\int_{-\infty}^{\infty} f_1(\tau)f_2(kT - \tau)d\tau$	$F_3(z)$	Continuous convolution does <i>not</i> correspond to product of z -transforms

B.2 TABLE OF z -TRANSFORMS

$\mathcal{F}(s)$ is the Laplace transform of $f(t)$ and $F(z)$ is the z -transform of $f(nT)$. Unless otherwise noted, $f(t) = 0$, $t < 0$ and the region of convergence of $F(z)$ is outside a circle $r < |z|$ such that all poles of $F(z)$ are inside r .

Table B.2

Number	$\mathcal{F}(s)$	$f(nT)$	$F(z)$
1	—	$1, n = 0; 0 n \neq 0$	1
2	—	$1, n = k; 0 n \neq k$	z^{-k}
3	$\frac{1}{s}$	$1(nT)$	$\frac{z}{z-1}$
4	$\frac{1}{s^2}$	nT	$\frac{Tz}{(z-1)^2}$
5	$\frac{1}{s^3}$	$\frac{1}{2!}(nT)^2$	$\frac{T^2 z(z+1)}{2(z-1)^3}$
6	$\frac{1}{s^4}$	$\frac{1}{3!}(nT)^3$	$\frac{T^3 z(z^2+4z+1)}{6(z-1)^4}$
7	$\frac{1}{s^m}$	$\lim_{a \rightarrow 0} \frac{(-1)^{m-1}}{(m-1)!} \frac{\partial^{m-1}}{\partial a^{m-1}} e^{-anT}$	$\lim_{a \rightarrow 0} \frac{(-1)^{m-1}}{(m-1)!} \frac{\partial^{m-1}}{\partial a^{m-1}} \frac{z}{z - e^{-aT}}$
8	$\frac{1}{s+a}$	e^{-anT}	$\frac{z}{z - e^{-aT}}$
9	$\frac{1}{(s+a)^2}$	nTe^{-anT}	$\frac{Tze^{-aT}}{(z - e^{-aT})^2}$
10	$\frac{1}{(s+a)^3}$	$\frac{1}{2}(nT)^2 e^{-anT}$	$\frac{T^2 e^{-aT} z(z + e^{-aT})}{2(z - e^{-aT})^3}$
11	$\frac{1}{(s+a)^m}$	$\frac{(-1)^{m-1}}{(m-1)!} \frac{\partial^{m-1}}{\partial a^{m-1}} (e^{-anT})$	$\frac{(-1)^{m-1}}{(m-1)!} \frac{\partial^{m-1}}{\partial a^{m-1}} \frac{z}{z - e^{-aT}}$
12	$\frac{a}{s(s+a)}$	$1 - e^{-anT}$	$\frac{z(1 - e^{-aT})}{(z-1)(z - e^{-aT})}$

Number	$\mathcal{F}(s)$	$f(nT)$	$F(z)$
13	$\frac{a}{s^2(s+a)}$	$\frac{1}{a}(anT - 1 + e^{-anT})$	$\frac{z[(aT - 1 + e^{-aT})z + (1 - e^{-aT} - aTe^{-aT})]}{a(z-1)^2(z - e^{-aT})}$
14	$\frac{b-a}{(s+a)(s+b)}$	$(e^{-anT} - e^{-bnT})$	$\frac{(e^{-aT} - e^{-bT})z}{(z - e^{-aT})(z - e^{-bT})}$
15	$\frac{s}{(s+a)^2}$	$(1 - anT)e^{-anT}$	$\frac{z[z - e^{-aT}(1 + aT)]}{(z - e^{-aT})^2}$
16	$\frac{a^2}{s(s+a)^2}$	$1 - e^{-anT}(1 + anT)$	$\frac{z[z(1 - e^{-aT} - aTe^{-aT}) + e^{-2aT} - e^{-aT} + aTe^{-aT}]}{(z-1)(z - e^{-aT})^2}$
17	$\frac{(b-a)s}{(s+a)(s+b)}$	$be^{-bnT} - ae^{-anT}$	$\frac{z[z(b-a) - (be^{-aT} - ae^{-bT})]}{(z - e^{-aT})(z - e^{-bT})}$
18	$\frac{a}{s^2 + a^2}$	$\sin anT$	$\frac{z \sin aT}{z^2 - (2 \cos aT)z + 1}$
19	$\frac{s}{s^2 + a^2}$	$\cos anT$	$\frac{z(z - \cos aT)}{z^2 - (2 \cos aT)z + 1}$
20	$\frac{s+a}{(s+a)^2 + b^2}$	$e^{-anT} \cos bnT$	$\frac{z(z - e^{-aT} \cos bT)}{z^2 - 2e^{-aT}(\cos bT)z + e^{-2aT}}$
21	$\frac{b}{(s+a)^2 + b^2}$	$e^{-anT} \sin bnT$	$\frac{ze^{-aT} \sin bT}{z^2 - 2e^{-aT}(\cos bT)z + e^{-2aT}}$
22	$\frac{a^2 + b^2}{s((s+a)^2 + b^2)}$	$1 - e^{-anT} \left(\cos bnT + \frac{a}{b} \sin bnT \right)$	$\frac{z(Az + B)}{(z-1)(z^2 - 2e^{-aT}(\cos bT)z + e^{-2aT})}$ $A = 1 - e^{-aT} \cos bT - \frac{a}{b} e^{-aT} \sin bT$ $B = e^{-2aT} + \frac{a}{b} e^{-aT} \sin bT - e^{-aT} \cos bT$

APPENDIX C

A Few Results from Matrix Analysis

Although we assume the reader has some acquaintance with linear equations and determinants, there are a few results of a more advanced character that even elementary control-system theory requires, and these are collected here for reference in the text. For further study, a good choice is Strang (1976).

C.1 DETERMINANTS AND THE MATRIX INVERSE

The determinant of a product of two square matrices is the product of their determinants:

$$\det \mathbf{AB} = \det \mathbf{A} \det \mathbf{B}. \quad (\text{C.1})$$

If a matrix is diagonal, then the determinant is the product of the elements on the diagonal.

If the matrix is partitioned with square elements on the main diagonal, then an extension of this result applies, namely,

$$\det \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{B} & \mathbf{C} \end{bmatrix} = \det \mathbf{A} \det \mathbf{C} \quad \text{if } \mathbf{A} \text{ and } \mathbf{C} \text{ are square matrices.} \quad (\text{C.2})$$

Suppose \mathbf{A} is a matrix of dimensions $m \times n$ and \mathbf{B} is of dimension $n \times m$. Let \mathbf{I}_m and \mathbf{I}_n be the identity matrices of size $m \times m$ and $n \times n$, respectively. Then

$$\det [\mathbf{I}_n + \mathbf{BA}] = \det [\mathbf{I}_m + \mathbf{AB}]. \quad (\text{C.3})$$

To show this result, we consider the determinant of the matrix product

$$\det \begin{bmatrix} \mathbf{I}_m & 0 \\ \mathbf{B} & \mathbf{I}_n \end{bmatrix} \begin{bmatrix} \mathbf{I}_m & \mathbf{A} \\ -\mathbf{B} & \mathbf{I}_n \end{bmatrix} = \det \begin{bmatrix} \mathbf{I}_m & \mathbf{A} \\ 0 & \mathbf{I}_n + \mathbf{BA} \end{bmatrix} = \det [\mathbf{I}_n + \mathbf{BA}].$$

But this is also equal to

$$\det \begin{bmatrix} \mathbf{I}_m & -\mathbf{A} \\ 0 & \mathbf{I}_n \end{bmatrix} \begin{bmatrix} \mathbf{I}_m & \mathbf{A} \\ -\mathbf{B} & \mathbf{I}_n \end{bmatrix} = \det \begin{bmatrix} \mathbf{I}_m + \mathbf{AB} & 0 \\ -\mathbf{B} & \mathbf{I}_n \end{bmatrix} = \det [\mathbf{I}_m + \mathbf{AB}],$$

and therefore these two determinants are equal to each other, which is (C.3).

If the determinant of a matrix \mathbf{A} is not zero, then we can define a related matrix \mathbf{A}^{-1} , called " \mathbf{A} inverse," which has the property that

$$\mathbf{AA}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}. \quad (\text{C.4})$$

According to property (C.1) we have

$$\det \mathbf{AA}^{-1} = \det \mathbf{A} \cdot \det \mathbf{A}^{-1} = 1,$$

or

$$\det \mathbf{A}^{-1} = \frac{1}{\det \mathbf{A}}.$$

It can be shown that there is an $n \times n$ matrix called the *adjugate* of \mathbf{A} with elements composed of sums of products of the elements of \mathbf{A}^1 and having the property that

$$\mathbf{A} \cdot \text{adj } \mathbf{A} = \det \mathbf{A} \cdot \mathbf{I}. \quad (\text{C.5})$$

Thus, if the determinant of \mathbf{A} is not zero, the inverse of \mathbf{A} is given by

$$\mathbf{A}^{-1} = \frac{\text{adj } \mathbf{A}}{\det \mathbf{A}}.$$

A famous and useful formula for the inverse of a combination of matrices has come to be called the *matrix inversion lemma* in the control literature. It arises in the development of recursive algorithms for estimation, as found

¹If \mathbf{A}^{ij} is the $n-1 \times n-1$ matrix (minor) found by deleting row i and column j from \mathbf{A} , then the entry in row i and column j of the $\text{adj } \mathbf{A}$ is $(-1)^{i+j} \det \mathbf{A}^{ji}$.

in Chapter 8. The formula is as follows: If $\det \mathbf{A}$, $\det \mathbf{C}$, and $\det (\mathbf{A} + \mathbf{BCD})$ are different from zero, then we have the matrix inversion lemma:

$$(\mathbf{A} + \mathbf{BCD})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{C}^{-1} + \mathbf{DA}^{-1}\mathbf{B})^{-1}\mathbf{DA}^{-1}. \quad (\text{C.6})$$

The truth of (C.6) is readily confirmed if we multiply both sides by $\mathbf{A} + \mathbf{BCD}$ to obtain

$$\begin{aligned} \mathbf{I} &= \mathbf{I} + \mathbf{BCDA}^{-1} - \mathbf{B}(\mathbf{C}^{-1} + \mathbf{DA}^{-1}\mathbf{B})^{-1}\mathbf{DA}^{-1} \\ &\quad - \mathbf{BCDA}^{-1}\mathbf{B}(\mathbf{C}^{-1} + \mathbf{DA}^{-1}\mathbf{B})^{-1}\mathbf{DA}^{-1} \\ &= \mathbf{I} + \mathbf{BCDA}^{-1} - [\mathbf{B} + \mathbf{BCDA}^{-1}\mathbf{B}][\mathbf{C}^{-1} + \mathbf{DA}^{-1}\mathbf{B})^{-1}\mathbf{DA}^{-1}. \end{aligned}$$

If we subtract \mathbf{I} from both sides and factor \mathbf{BC} from the left on the third term, we find

$$0 = \mathbf{BCDA}^{-1} - \mathbf{BC}[\mathbf{C}^{-1} + \mathbf{DA}^{-1}\mathbf{B}][\mathbf{C}^{-1} + \mathbf{DA}^{-1}\mathbf{B})^{-1}\mathbf{DA}^{-1},$$

which is

$$0 = 0 \text{ which was to be demonstrated.}$$

C.2 EIGENVALUES AND EIGENVECTORS

We consider the discrete dynamic system

$$\mathbf{x}_{k+1} = \Phi \mathbf{x}_k, \quad (\text{C.7})$$

where, for purposes of illustration, we will let

$$\Phi = \begin{bmatrix} \frac{5}{6} & -\frac{1}{6} \\ 1 & 0 \end{bmatrix}. \quad (\text{C.8})$$

If we assume that it is possible for this system to have a motion given by a geometric series such as z^k , we can assume that there is a vector \mathbf{v} so that \mathbf{x}_k can be written

$$\mathbf{x}_k = \mathbf{v}z^k. \quad (\text{C.9})$$

Substituting (C.9) into (C.7), we must find the vector \mathbf{v} and the number z such that

$$\mathbf{v}z^{k+1} = \Phi \mathbf{v}z^k,$$

or, multiplying by z^{-k} yields

$$\mathbf{v}z = \Phi \mathbf{v}. \quad (\text{C.10})$$

If we collect both the terms of (C.10) on the left, we find

$$(z\mathbf{I} - \Phi)\mathbf{v} = 0. \quad (\text{C.11})$$

These linear equations have a solution for a nontrivial \mathbf{v} if and only if the determinant of the coefficient matrix is zero. This determinant is a polynomial of degree n in z (Φ is an $n \times n$ matrix) called the *characteristic polynomial* of Φ , and values of z for which the characteristic polynomial is zero are roots of the characteristic equation and are called *eigenvalues* of Φ . For example, for the matrix given in (C.8) the characteristic polynomial is

$$\det \left\{ \begin{bmatrix} z & 0 \\ 0 & z \end{bmatrix} - \begin{bmatrix} \frac{5}{6} & -\frac{1}{6} \\ 1 & 0 \end{bmatrix} \right\}.$$

Adding the two matrices, we find

$$\det \left\{ \begin{bmatrix} z - \frac{5}{6} & +\frac{1}{6} \\ -1 & z \end{bmatrix} \right\},$$

which can be evaluated to give

$$z(z - \frac{5}{6}) + \frac{1}{6} = (z - \frac{1}{2})(z - \frac{1}{3}). \quad (\text{C.12})$$

Thus the characteristic roots of this Φ are $\frac{1}{2}$ and $\frac{1}{3}$. Associated with these characteristic roots are solutions to (C.11) for vectors \mathbf{v} , called the *characteristic* or *eigenvectors*. If we let $z = \frac{1}{2}$, then (C.11) requires

$$\left\{ \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{bmatrix} - \begin{bmatrix} \frac{5}{6} & -\frac{1}{6} \\ 1 & 0 \end{bmatrix} \right\} \begin{bmatrix} v_{11} \\ v_{21} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}. \quad (\text{C.13})$$

Adding the matrices, we find that these equations become

$$\begin{bmatrix} -\frac{1}{3} & \frac{1}{6} \\ -1 & \frac{1}{2} \end{bmatrix} \begin{bmatrix} v_{11} \\ v_{21} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad (\text{C.14})$$

Equations (C.14) are satisfied by any v_{11} and v_{21} such that

$$v_{21} = 2v_{11},$$

from which we conclude that the eigenvector corresponding to $z_1 = \frac{1}{2}$ is given by

$$\mathbf{v}_1 = \begin{bmatrix} a \\ 2a \end{bmatrix}. \quad (\text{C.15})$$

We can arbitrarily select the scale factor a in (C.15). Some prefer to make the length² of eigenvectors equal to one. Here we make the largest component of \mathbf{v} have unit magnitude. Thus the scaled \mathbf{v}_1 is

$$\mathbf{v}_1 = \begin{bmatrix} \frac{1}{2} \\ 1 \end{bmatrix}. \quad (\text{C.16})$$

In similar fashion, the eigenvector \mathbf{v}_2 associated with $z_2 = \frac{1}{3}$ can be computed to be

$$\mathbf{v}_2 = \begin{bmatrix} \frac{1}{3} \\ 1 \end{bmatrix}.$$

Note that even if all elements of Φ are real, it is possible for characteristic values and characteristic vectors to be complex.

C.3 SIMILARITY TRANSFORMATIONS

If we make a change of variables in (C.7) according to $\mathbf{x} = \mathbf{T}\xi$, where \mathbf{T} is an $n \times n$ matrix, then we start with the equations

$$\mathbf{x}_{k+1} = \Phi \mathbf{x}_k,$$

and, substituting for \mathbf{x} , we have

$$\mathbf{T}\xi_{k+1} = \Phi \mathbf{T}\xi_k.$$

Then, if we multiply on the left by \mathbf{T}^{-1} , we get the equation in ξ ,

$$\xi_{k+1} = \mathbf{T}^{-1}\Phi \mathbf{T}\xi_k. \quad (\text{C.17})$$

²Usually we define the length of a vector as the square root of the sum of squares of its components or, if $\|\mathbf{v}\|$ is the symbol for length, then $\|\mathbf{v}\|^2 = \mathbf{v}^T \mathbf{v}$. If \mathbf{v} is complex, as will happen if z_i is complex, then we must take a conjugate, and we define $\|\mathbf{v}\|^2 = (\mathbf{v}^*)^T \mathbf{v}$, where \mathbf{v}^* is the complex conjugate of \mathbf{v} .

and the covariance matrix,

$$\mathcal{E}\{(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^T\} = \mathbf{R}_x.$$

Like the scalar normal density, the multivariable law is described entirely by the two parameters $\boldsymbol{\mu}$ and \mathbf{R} , the difference being that the multivariable case is described by matrix parameters rather than scalar parameters. In (D.19) we require the inverse of \mathbf{R}_x and have thus implicitly assumed that this covariance matrix is nonsingular. [See Parzen (1962) for a discussion of the case when \mathbf{R}_x is singular.]

D.4 STOCHASTIC PROCESSES

In a study of dynamic systems, it is natural to have random variables that evolve in time much as the states and control inputs evolve. However, with random time variables it is not possible to compute z -transforms in the usual way; and furthermore, because specific values of the variables have little value, we need formulas to describe how the means and covariances evolve in time. A random variable that evolves in time is called a *stochastic process*, and here we consider only discrete time.

Suppose we deal first with a stochastic process $w(n)$, where w is a scalar distributed according to the density $f_w(\xi; N)$. Note that the density function depends on the time of occurrence of the random variable. If a variable has statistical properties (such as f_w) that are independent of the origin of time, then we say the process is *stationary*. Considering values of the process at distinct times, we have separate random variables, and we define the covariance of the process w as

$$R_w(j, k) = \mathcal{E}(w(j) - \bar{w}(j))(w(k) - \bar{w}(k)). \quad (\text{D.20})$$

If the process is stationary, then the covariance in (D.20) depends only on the magnitude of the difference in observation times, $k - j$, and we often will write $R_w(j, k) = R_w(k - j)$ and drop the second argument. Because a stochastic process is both random and time dependent, we can imagine averages that are computed over the time variable as well as by the expectation. For example, for a stationary process $w(n)$ we can define the mean as

$$\bar{w}(k) = \lim_{N \rightarrow \infty} \frac{1}{2N + 1} \sum_{n=-N}^N w(n + k), \quad (\text{D.21})$$

and the second-order mean or autocorrelation

$$\begin{aligned} & (\widetilde{w(j)} - \widetilde{w})(\widetilde{w(k)} - \widetilde{w(k)}) \\ &= \lim_{N \rightarrow \infty} \frac{1}{2N+1} \sum_{n=-N}^N \{(w(n+j) - \widetilde{w(j)})(w(n+k) - \widetilde{w(k)})\}. \end{aligned} \quad (\text{D.22})$$

For a stationary process, the time average in (D.21) is usually equal to the distribution average, and likewise the second-order average in (D.22) is the same as the covariance in (D.20). Processes for which time averages give the same limits as distribution averages are called *ergodic*.

A very useful aid to understanding the properties of stationary stochastic processes is found by considering the response of a linear stationary system to a stationary input process. Suppose we let the input be w , a stationary scalar process with zero mean and covariance $R_w(j)$, and suppose we take the output to be $y(k)$. We let the unit-pulse response from w to y be $h(j)$. Thus from standard analysis (see Chapter 2), we have

$$y(j) = \sum_{k=-\infty}^{\infty} h(k)w(j-k), \quad (\text{D.23})$$

and the covariance of $y(j)$ with $y(j+\ell)$ is

$$\begin{aligned} R_y(\ell) &= \mathcal{E}y(j+\ell)y(j) \\ &= \mathcal{E} \left\{ \sum_{k=-\infty}^{\infty} h(k)w(j+\ell-k) \right\} \left\{ \sum_{n=-\infty}^{\infty} h(n)w(j-n) \right\}. \end{aligned} \quad (\text{D.24})$$

Because the system unit-pulse response, $h(k)$, is not random, both $h(k)$ and $h(n)$ can be removed from the integral implied by the \mathcal{E} operation, with the result

$$R_y(\ell) = \sum_{k=-\infty}^{\infty} h(k) \sum_{n=-\infty}^{\infty} h(n) \mathcal{E}w(j+\ell-k)w(j-n). \quad (\text{D.25})$$

The expectation in (D.25) is now recognized as $R_w(\ell-k+n)$, and substituting this expression in (D.25), we find

$$R_y(\ell) = \sum_{k=-\infty}^{\infty} h(k) \sum_{n=-\infty}^{\infty} h(n) R_w(\ell-k+n). \quad (\text{D.26})$$

7

System Time Response Characteristics

In this chapter we investigate the time response of a sampled data system and compare it with the response of a similar continuous system. In addition, the mapping between the s -domain and the z -domain is examined, the important time response characteristics of continuous systems are revised and their equivalents in the discrete domain are discussed.

7.1 TIME RESPONSE COMPARISON

An example closed-loop discrete-time system with a zero-order hold is shown in Figure 7.1(a). The continuous-time equivalent of this system is also shown in Figure 7.1(b), where the sampler (A/D converter) and the zero-order hold (D/A converter) have been removed. We shall now derive equations for the step responses of both systems and then plot and compare them.

As described in Chapter 6, the transfer function of the above discrete-time system is given by

$$\frac{y(z)}{r(z)} = \frac{G(z)}{1 + G(z)}, \quad (7.1)$$

where

$$r(z) = \frac{z}{z - 1} \quad (7.2)$$

and the z -transform of the plant is given by

$$G(s) = \frac{1 - e^{-sT}}{s^2(s + 1)}.$$

Expanding by means of partial fractions, we obtain

$$G(s) = (1 - e^{-sT}) \left(\frac{1}{s^2} - \frac{1}{s} + \frac{1}{s + 1} \right)$$

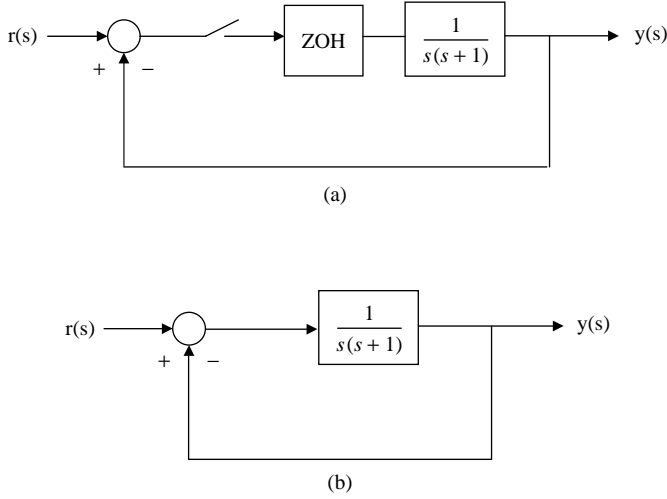


Figure 7.1 (a) Discrete system and (b) its continuous-time equivalent

and the z -transform is

$$G(z) = (1 - z^{-1})Z \left\{ \frac{1}{s^2} - \frac{1}{s} + \frac{1}{s+1} \right\}.$$

From z -transform tables we obtain

$$G(z) = (1 - z^{-1}) \left[\frac{Tz}{(z-1)^2} - \frac{z}{z-1} + \frac{z}{z-e^{-T}} \right].$$

Setting $T = 1$ s and simplifying gives

$$G(z) = \frac{0.368z + 0.264}{z^2 - 1.368z + 0.368}.$$

Substituting into (7.1), we obtain the transfer function

$$\frac{y(z)}{r(z)} = \frac{G(z)}{1 + G(z)} = \frac{0.368z + 0.264}{z^2 - z + 0.632},$$

and then using (7.2) gives the output

$$y(z) = \frac{z(0.368z + 0.264)}{(z-1)(z^2 - z + 0.632)}.$$

The inverse z -transform can be found by long division: the first several terms are

$$y(z) = 0.368z^{-1} + z^{-2} + 1.4z^{-3} + 1.4z^{-4} + 1.15z^{-5} + 0.9z^{-6} + 0.8z^{-7} + 0.87z^{-8} \\ + 0.99z^{-9} + \dots$$

and the time response is given by

$$y(nT) = 0.368\delta(t-1) + \delta(t-2) + 1.4\delta(t-3) + 1.4\delta(t-4) + 1.15\delta(t-5) \\ + 0.9\delta(t-6) + 0.8\delta(t-7) + 0.87\delta(t-8) + \dots$$

From Figure 7.1(b), the equivalent continuous-time system transfer function is

$$\frac{y(s)}{r(s)} = \frac{G(s)}{1 + G(s)} = \frac{1/(s(s+1))}{1 + (1/(s(s+1)))} = \frac{1}{s^2 + s + 1}.$$

Since $r(s) = 1/s$, the output becomes

$$y(s) = \frac{1}{s(s^2 + s + 1)}.$$

To find the inverse Laplace transform we can write

$$y(s) = \frac{1}{s} - \frac{s+1}{s^2 + s + 1} = \frac{1}{s} - \frac{s+0.5}{(s+0.5)^2 - 0.5^2} - \frac{0.5}{(s+0.5)^2 - 0.5^2}.$$

From inverse Laplace transform tables we find that the time response is

$$y(t) = 1 - e^{-0.5t} (\cos 0.5t + 0.577 \sin 0.5t).$$

Figure 7.2 shows the time responses of both the discrete-time system and its continuous-time equivalent. The response of the discrete-time system is accurate only at the sampling instants. As shown in the figure, the sampling process has a destabilizing effect on the system.

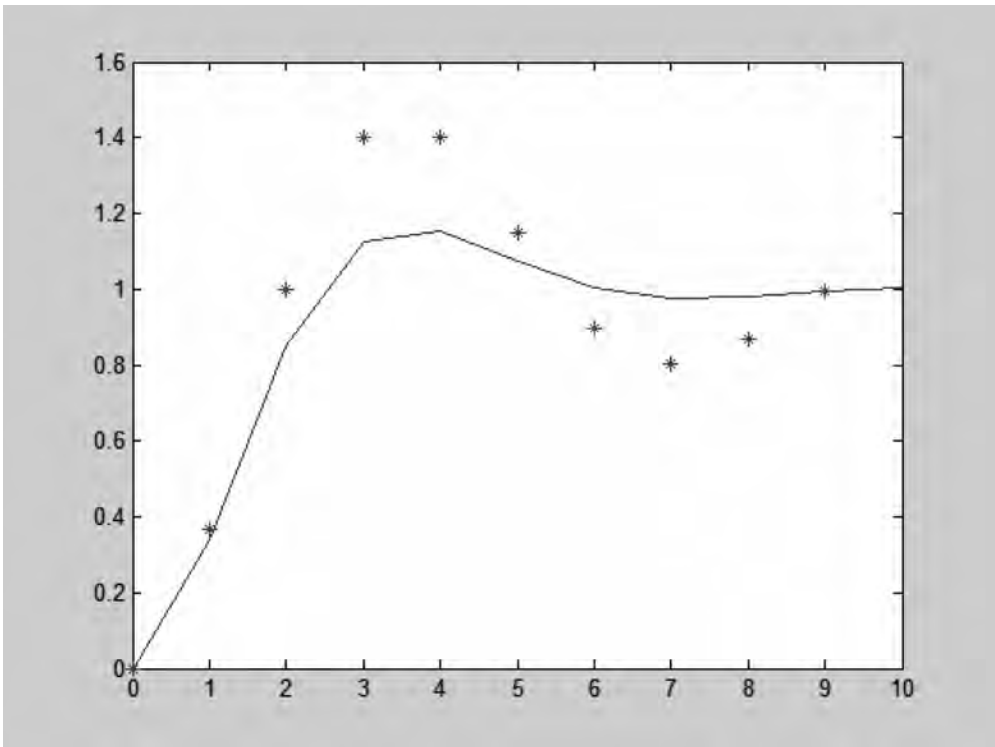


Figure 7.2 Step response of the system shown in Figure 7.1

7.2 TIME DOMAIN SPECIFICATIONS

The performance of a control system is usually measured in terms of its response to a step input. The step input is used because it is easy to generate and gives the system a nonzero steady-state condition, which can be measured.

Most commonly used time domain performance measures refer to a second-order system with the transfer function:

$$\frac{y(s)}{r(s)} = \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2},$$

where ω_n is the undamped natural frequency of the system and ζ is the damping ratio of the system.

When a second-order system is excited with a unit step input, the typical output response is as shown in Figure 7.3. Based on this figure, the following performance parameters are usually defined: maximum overshoot; peak time; rise time; settling time; and steady-state error.

The maximum overshoot, M_p , is the peak value of the response curve measured from unity. This parameter is usually quoted as a percentage. The amount of overshoot depends on the damping ratio and directly indicates the relative stability of the system.

The peak time, T_p , is defined as the time required for the response to reach the first peak of the overshoot. The system is more responsive when the peak time is smaller, but this gives rise to a higher overshoot.

The rise time, T_r , is the time required for the response to go from 0 % to 100 % of its final value. It is a measure of the responsiveness of a system, and smaller rise times make the system more responsive.

The settling time, T_s , is the time required for the response curve to reach and stay within a range about the final value. A value of 2–5 % is usually used in performance specifications.

The steady-state error, E_{ss} , is the error between the system response and the reference input value (unity) when the system reaches its steady-state value. A small steady-state error is a requirement in most control systems. In some control systems, such as position control, it is one of the requirements to have no steady-state error.

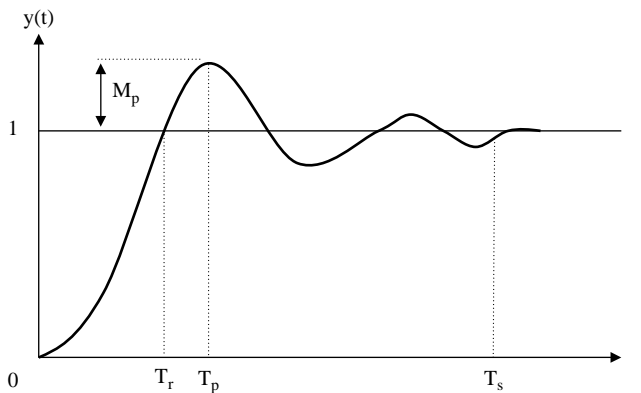


Figure 7.3 Second-order system unit step response

Having introduced the parameters, we are now in a position to give formulae for them (readers who are interested in the derivation of these formulae should refer to books on control theory). The maximum overshoot occurs at peak time ($t = T_p$) and is given by

$$M_p = e^{-(\zeta\pi/\sqrt{1-\zeta^2})},$$

i.e. overshoot is directly related to the system damping ratio – the lower the damping ratio, the higher the overshoot. Figure 7.4 shows the variation of the overshoot (expressed as a percentage) with the damping ratio.

The peak time is obtained by differentiating the output response with respect to time, letting this equal zero. It is given by

$$T_p = \frac{\pi}{\omega_d},$$

where

$$\omega_d = \omega_n^2 \sqrt{1 - \zeta^2}$$

is the damped natural frequency.

The rise time is obtained by setting the output response to 1 and finding the time. It is given by

$$T_r = \frac{\pi - \beta}{\omega_d},$$

where

$$\beta = \tan^{-1} \frac{\omega_d}{\zeta \omega_n}.$$

The settling time is usually specified for a 2 % or 5 % tolerance band, and is given by

$$T_s = \frac{4}{\zeta \omega_n} \quad (\text{for 2\% settling time}),$$

$$T_s = \frac{3}{\zeta \omega_n} \quad (\text{for 5\% settling time}).$$

The steady-state error can be found by using the final value theorem, i.e. if the Laplace transform of the output response is $y(s)$, then the final value (steady-state value) is given by

$$\lim_{s \rightarrow 0} sy(s),$$

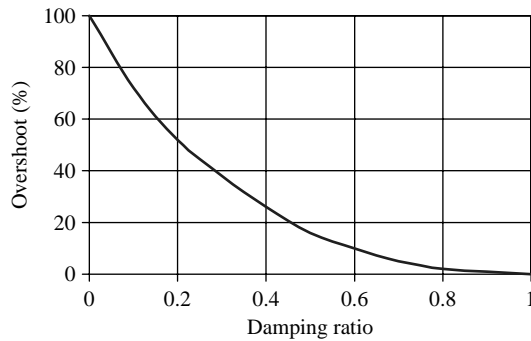


Figure 7.4 Variation of overshoot with damping ratio

and the steady-state error when a unit step input is applied can be found from

$$E_{ss} = 1 - \lim_{s \rightarrow 0} s y(s).$$

Example 7.1

Determine the performance parameters of the system given in Section 7.1 with closed-loop transfer function

$$\frac{y(s)}{r(s)} = \frac{1}{s^2 + s + 1}.$$

Solution

Comparing this system with the standard second-order system transfer function

$$\frac{y(s)}{r(s)} = \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2},$$

we find that $\zeta = 0.5$ and $\omega_n = 1$ rad/s. Thus, the damped natural frequency is

$$\omega_d = \omega_n^2 \sqrt{1 - \zeta^2} = 0.866 \text{ rad/s}.$$

The peak overshoot is

$$M_p = e^{-(\zeta\pi/\sqrt{1-\zeta^2})} = 0.16$$

or 16 %. The peak time is

$$T_p = \frac{\pi}{\omega_d} = 3.627 \text{ s}$$

The rise time is

$$T_r = \frac{\pi - \beta}{\omega_n};$$

since

$$\beta = \tan^{-1} \frac{\omega_d}{\zeta\omega_n} = 1.047,$$

we have

$$T_r = \frac{\pi - \beta}{\omega_n} = \frac{\pi - 1.047}{1} = 2.094 \text{ s}$$

The settling time (2 %) is

$$T_s = \frac{4}{\zeta\omega_n} = 8 \text{ s},$$

and the settling time (5 %) is

$$T_s = \frac{3}{\zeta\omega_n} = 6 \text{ s}.$$

Finally, the steady state error is

$$E_{ss} = 1 - \lim_{s \rightarrow 0} s y(s) = 1 - \lim_{s \rightarrow 0} s \frac{1}{s(s^2 + s + 1)} = 0.$$

7.3 MAPPING THE s -PLANE INTO THE z -PLANE

The pole locations of a closed-loop continuous-time system in the s -plane determine the behaviour and stability of the system, and we can shape the response of a system by positioning its poles in the s -plane. It is desirable to do the same for the sampled data systems. This section describes the relationship between the s -plane and the z -plane and analyses the behaviour of a system when the closed-loop poles are placed in the z -plane.

First of all, consider the mapping of the left-hand side of the s -plane into the z -plane. Let $s = \sigma + j\omega$ describe a point in the s -plane. Then, along the $j\omega$ axis,

$$z = e^{sT} = e^{\sigma T} e^{j\omega T}.$$

But $\sigma = 0$ so we have

$$z = e^{j\omega T} = \cos \omega T + j \sin \omega T = 1 \angle \omega T.$$

Hence, the pole locations on the imaginary axis in the s -plane are mapped onto the unit circle in the z -plane. As ω changes along the imaginary axis in the s -plane, the angle of the poles on the unit circle in the z -plane changes.

If ω is kept constant and σ is increased in the left-hand s -plane, the pole locations in the z -plane move towards the origin, away from the unit circle. Similarly, if σ is decreased in the left-hand s -plane, the pole locations in the z -plane move away from the origin in the z -plane. Hence, the entire left-hand s -plane is mapped into the interior of the unit circle in the z -plane. Similarly, the right-hand s -plane is mapped into the exterior of the unit circle in the z -plane. As far as the system stability is concerned, a sampled data system will be stable if the closed-loop poles (or the zeros of the characteristic equation) lie within the unit circle. Figure 7.5 shows the mapping of the left-hand s -plane into the z -plane.

As shown in Figure 7.6, lines of constant σ in the s -plane are mapped into circles in the z -plane with radius $e^{\sigma T}$. If the line is on the left-hand side of the s -plane then the radius of the circle in the z -plane is less than 1. If on the other hand the line is on the right-hand side of the s -plane then the radius of the circle in the z -plane is greater than 1. Figure 7.7 shows the corresponding pole locations between the s -plane and the z -plane.

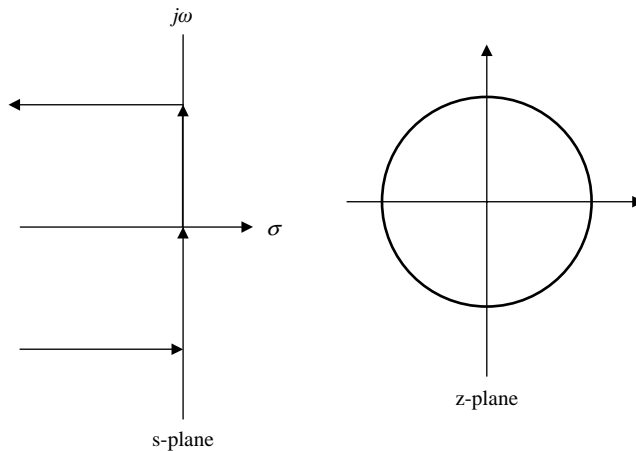


Figure 7.5 Mapping the left-hand s -plane into the z -plane

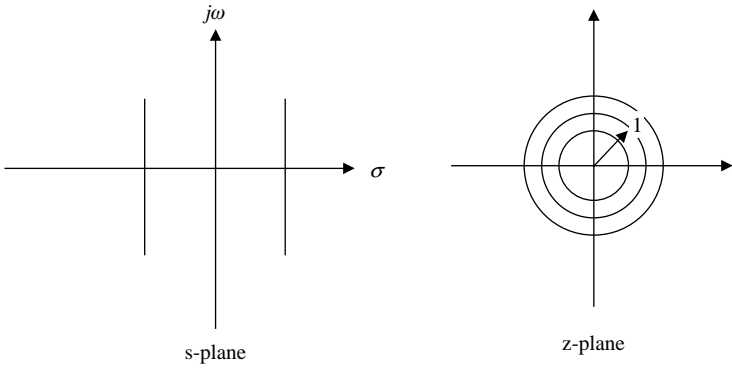


Figure 7.6 Mapping the lines of constant σ

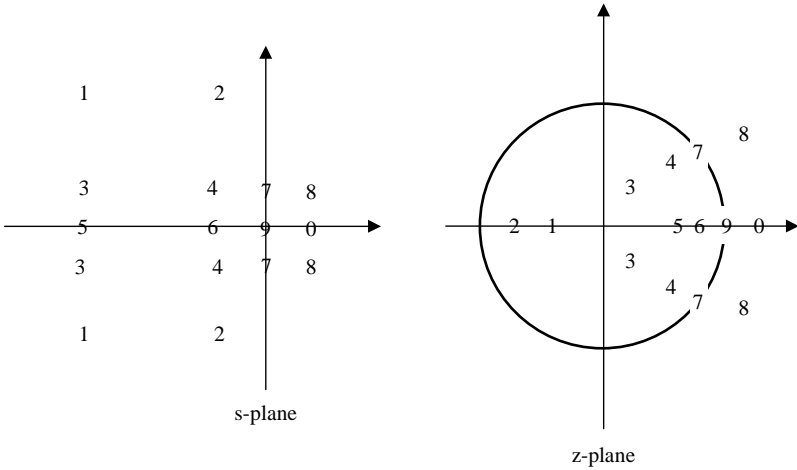


Figure 7.7 Poles in the s-plane and their corresponding z-plane locations

The time responses of a sampled data system based on its pole positions in the z-plane are shown in Figure 7.8. It is clear from this figure that the system is stable if all the closed-loop poles are within the unit circle.

7.4 DAMPING RATIO AND UNDAMPED NATURAL FREQUENCY IN THE z-PLANE

7.4.1 Damping Ratio

As shown in Figure 7.9(a), lines of constant damping ratio in the s-plane are lines where $\zeta = \cos \alpha$ for a given damping ratio. The locus in the z-plane can then be obtained by the substitution $z = e^{sT}$. Remembering that we are working in the third and fourth quadrants in

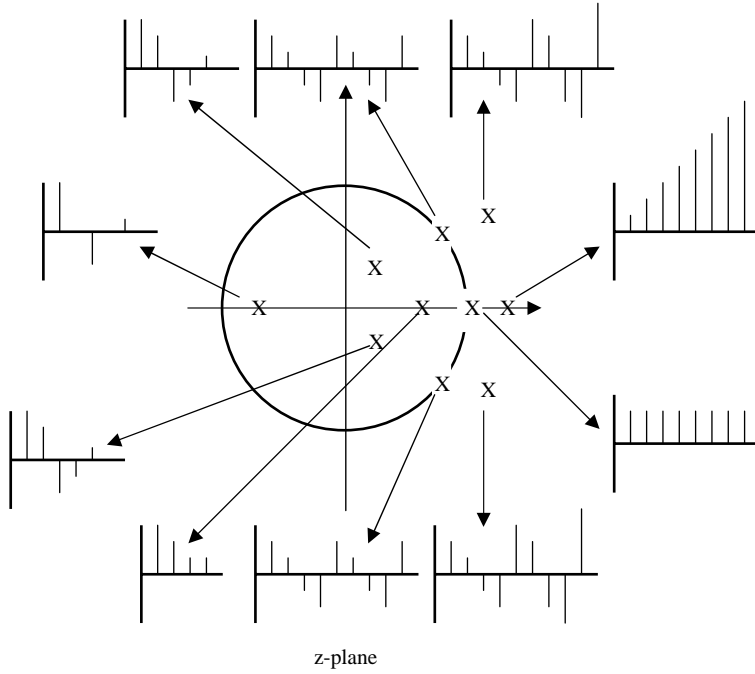


Figure 7.8 Time response of z-plane pole locations

the s -plane where s is negative, we get

$$z = e^{-\sigma\omega T} e^{j\omega T}. \quad (7.3)$$

Since, from Figure 7.9(a),

$$\sigma = \tan\left(\frac{\pi}{2} - \cos^{-1} \zeta\right), \quad (7.4)$$

substituting in (7.3) we have

$$z = \exp\left[-\omega T \tan\left(\frac{\pi}{2} - \cos^{-1} \zeta\right)\right] e^{j\omega T}. \quad (7.5)$$

Equation (7.5) describes a logarithmic spiral in the z -plane as shown in Figure 7.9(b). The spiral starts from $z = 1$ when $\omega = 0$. Figure 7.10 shows the lines of constant damping ratio in the z -plane for various values of ζ .

7.4.2 Undamped Natural Frequency

As shown in Figure 7.11, the locus of constant undamped natural frequency in the s -plane is a circle with radius ω_n . From this figure, we can write

$$\omega^2 + \sigma^2 = \omega_n^2 \quad \text{or} \quad \sigma = \sqrt{\omega_n^2 - \omega^2}. \quad (7.6)$$

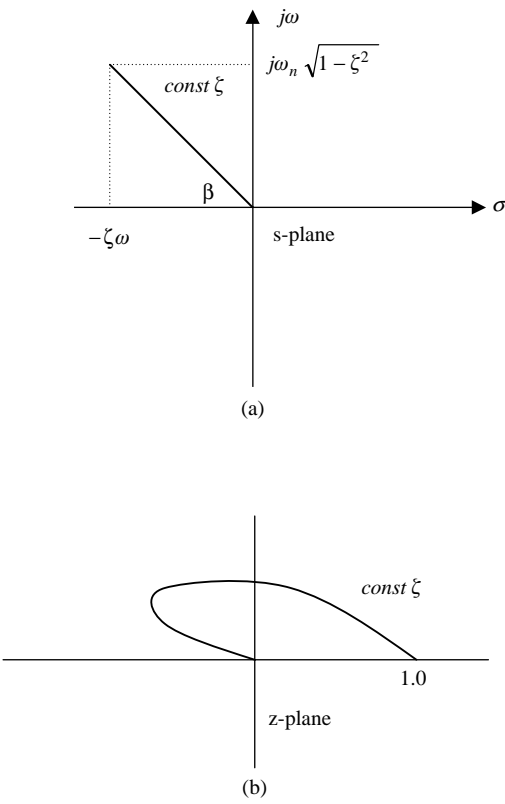


Figure 7.9 (a) Line of constant damping ratio in the s -plane, and (b) the corresponding locus in the z -plane

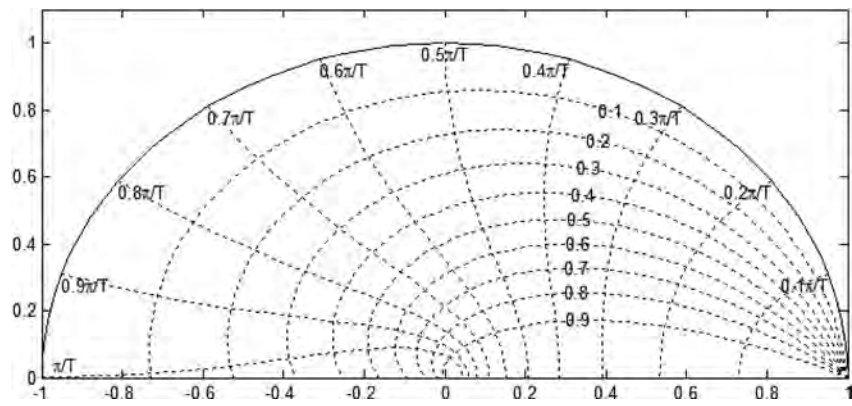


Figure 7.10 Lines of constant damping ratio for different ζ . The vertical lines are the lines of constant ω_n

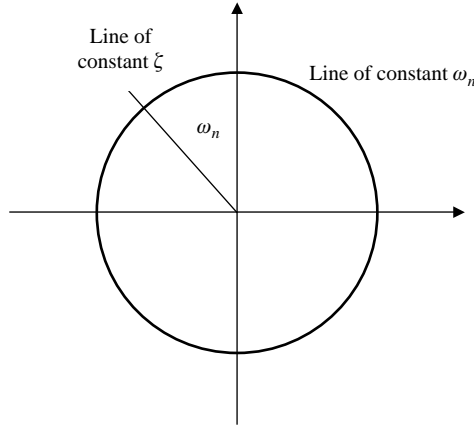


Figure 7.11 Locus of constant ω_n in the s -plane

Thus, remembering that s is negative, we have

$$z = e^{-sT} = e^{-\sigma T} e^{-j\omega T} = \exp \left[-T(\sqrt{\omega_n^2 - \omega^2}) \right] e^{-j\omega T} \quad (7.7)$$

The locus of constant ω_n in the z -plane is given by (7.7) and is shown in Figure 7.10 as the vertical lines. Notice that the curves are given for values of ω_n ranging from $\omega_n = \pi/10T$ to $\omega_n = \pi/T$.

Notice that the loci of constant damping ratio and the loci of undamped natural frequency are usually shown on the same graph.

7.5 DAMPING RATIO AND UNDAMPED NATURAL FREQUENCY USING FORMULAE

In Section 7.4 above we saw how to find the damping ratio and the undamped natural frequency of a system using a graphical technique. Here, we will derive equations for calculating the damping ratio and the undamped natural frequency.

The damping ratio and the natural frequency of a system in the z -plane can be determined if we first of all consider a second-order system in the s -plane:

$$G(s) = \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2}. \quad (7.8)$$

The poles of this system are at

$$s_{1,2} = -\zeta\omega_n \pm j\omega_n\sqrt{1 - \zeta^2}. \quad (7.9)$$

We can now find the equivalent z -plane poles by making the substitution $z = e^{sT}$, i.e.

$$z = e^{sT} = e^{-\zeta\omega_n T} \angle \pm \omega_n T \sqrt{1 - \zeta^2}, \quad (7.10)$$

which we can write as

$$z = r \angle \pm \theta, \quad (7.11)$$

where

$$r = e^{-\zeta \omega_n T} \quad \text{or} \quad \zeta \omega_n T = -\ln r \quad (7.12)$$

and

$$\theta = \omega_n T \sqrt{1 - \zeta^2}. \quad (7.13)$$

From (7.12) and (7.13) we obtain

$$\frac{\zeta}{\sqrt{1 - \zeta^2}} = \frac{-\ln r}{\theta}$$

or

$$\zeta = \frac{-\ln r}{\sqrt{(\ln r)^2 + \theta^2}}, \quad (7.14)$$

and from (7.12) and (7.14) we obtain

$$\omega_n = \frac{1}{T} \sqrt{(\ln r)^2 + \theta^2}. \quad (7.15)$$

Example 7.2

Consider the system described in Section 7.1 with closed-loop transfer function

$$\frac{y(z)}{r(z)} = \frac{G(z)}{1 + G(z)} = \frac{0.368z + 0.264}{z^2 - z + 0.632}.$$

Find the damping ratio and the undamped natural frequency. Assume that $T = 1$ s.

Solution

We need to find the poles of the closed-loop transfer function. The system characteristic equation is $1 + G(z) = 0$,

i.e.

$$z^2 - z + 0.632 = (z - 0.5 - j0.618)(z - 0.5 + j0.618) = 0,$$

which can be written in polar form as

$$z_{1,2} = 0.5 \pm j0.618 = 0.795 \angle \pm 0.890 = r \angle \pm \theta$$

(see (7.11)). The damping ratio is then calculated using (7.14) as

$$\zeta = \frac{-\ln r}{\sqrt{(\ln r)^2 + \theta^2}} = \frac{-\ln 0.795}{\sqrt{(\ln 0.795)^2 + 0.890^2}} = 0.25,$$

and from (7.15) the undamped natural frequency is, taking $T = 1$,

$$\omega_n = \frac{1}{T} \sqrt{(\ln r)^2 + \theta^2} = \sqrt{(\ln 0.795)^2 + 0.890^2} = 0.92.$$

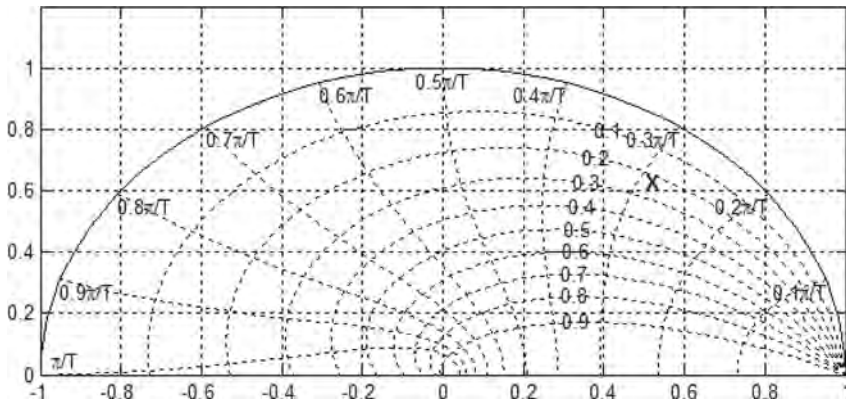


Figure 7.12 Finding ζ and ω_n graphically

Example 7.3

Find the damping ratio and the undamped natural frequency for Example 7.2 using the graphical method.

Solution

The characteristic equation of the system is found to be

$$z^2 - z + 0.632 = (z - 0.5 - j0.618)(z - 0.5 + j0.618) = 0$$

and the poles of the closed-loop system are at

$$z_{1,2} = 0.5 \pm j0.618.$$

Figure 7.12 shows the loci of the constant damping ratio and the loci of the undamped natural frequency with the poles of the closed-loop system marked with an \times on the graph. From the graph we can read the damping ratio as 0.25 and the undamped natural frequency as

$$\omega_n = \frac{0.29\pi}{T} = 0.91.$$

7.6 EXERCISES

- Find the damping ratio and the undamped natural frequency of the sampled data systems whose characteristic equations are given below
 - $z^2 - z + 2 = 0$
 - $z^2 - 1 = 0$
 - $z^2 - z + 1 = 0$
 - $z^2 - 0.81 = 0$

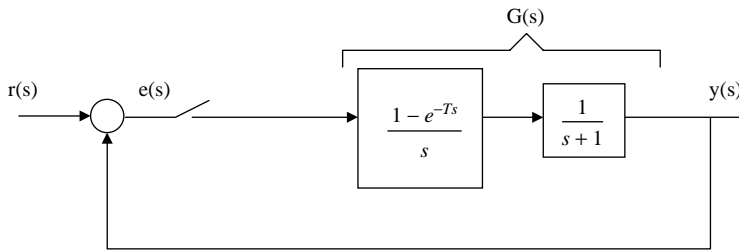


Figure 7.13 System for Exercise 2

2. Consider the closed-loop system of Figure 7.13. Assume that $T = 1$ s.
 - (a) Calculate the transfer function of the system.
 - (b) Calculate and plot the unit step response at the sampling instants.
 - (c) Calculate the damping factor and the undamped natural frequency of the system.
3. Consider the closed-loop system of Figure 7.13. Do not assume a value for T .
 - (a) Calculate the transfer function of the system.
 - (b) Calculate the damping factor and the undamped natural frequency of the system.
 - (c) What will be the steady state error if a unit step input is applied?
4. A unit step input is applied to the system in Figure 7.13. Calculate:
 - (a) the percentage overshoot;
 - (b) the peak time;
 - (c) the rise time;
 - (d) settling time to 5 %.
5. The closed-loop transfer functions of four sampled data systems are given below. Calculate the percentage overshoots and peak times.
 - (a) $G(z) = \frac{1}{z^2 + z + 2}$
 - (b) $G(z) = \frac{1}{z^2 + 2z + 1}$
 - (c) $G(z) = \frac{1}{z^2 - z + 1}$
 - (d) $G(z) = \frac{2}{z^2 + z + 4}$
6. The s -plane poles of a continuous-time system are at $s = -1$ and $s = -2$. Assuming $T = 1$ s, calculate the pole locations in the z -plane.
7. The s -plane poles of a continuous-time system are at $s_{1,2} = -0.5 \pm j0.9$. Assuming $T = 1$ s, calculate the pole locations in the z -plane. Calculate the damping ratio and the undamped natural frequency of the system using a graphical technique.

FURTHER READING

- [D'Azzo and Houpis, 1966] D'Azzo, J.J. and Houpis, C.H. Feedback Control System Analysis and Synthesis, 2nd edn., McGraw-Hill, New York, 1966.
- [Dorf, 1992] Dorf, R.C. Modern Control Systems, 6th edn. Addison-Wesley, Reading, MA, 1992.

in the fact that the control response is determined from

$$\frac{U(z)}{R(z)} = \frac{D}{1 + DG} = \frac{H(z)}{G(z)},$$

which for this example is

$$\frac{U(z)}{R(z)} = 13.06 \frac{z - 0.0793}{z^2 - 0.7859z + 0.3679} \frac{(z - 1)(z - 0.9048)}{z + 0.9672}.$$

There is a root at $z = -0.9672$!. This is the source of the oscillation in the control response, but it did not show up in the output response because it was exactly canceled by a zero. The control oscillation causes the “intersample ripple” in the output response, and the designer should be alert to this if poorly behaved roots arise in the control response. An actual prediction of the output intersample ripple based on linear analysis was not possible with the z -transform method described so far; rather, one would need to apply the “modified z -transform,” which is beyond the scope of this text. Alternatively, one can use a CAD simulation to find such oscillations quite easily, as was done here. To avoid this oscillation, we could introduce another term in $H(z)$, $b_3 z^{-3}$, and require that $H(z)$ be zero at $z = -0.9672$, so this zero of $G(z)$ is not canceled by $D(z)$. The result will be a simpler $D(z)$ with a slightly more complicated $H(z)$. However, rather than pursue this method further, we will wait until the more powerful method of pole assignment by state-variable analysis is developed in the next chapter, where computer algorithms are more readily provided.

5.8 PID CONTROL

Just as in continuous systems, there are three basic types of control: Proportional, Integral, and Derivative, hence the name, PID. In the design examples so far, we have been using the discrete equivalent of lead compensation, which is essentially a combination of proportional and derivative control. Let us now review these three controls as they pertain to a discrete implementation. The term PID is widely used because there are commercially available modules that have knobs for the user to turn that set the values of each of the three control types.

5.8.1 Proportional Control

A discrete implementation of proportional control is identical to continuous; that is, where the continuous is

$$u(t) = K_p e(t) \Rightarrow D(s) = K_p,$$

the discrete is

$$u(k) = K_p e(k) \Rightarrow \boxed{D(z) = K_p}$$

where $e(t)$ is the error signal as shown in Fig 5.2.

5.8.2 Derivative Control

For continuous systems, derivative or rate control has the form

$$u(t) = K_p T_D \dot{e}(t) \Rightarrow D(s) = K_p T_D s$$

where T_D is called the *derivative time*. Differentiation can be approximated in the discrete domain as the first difference, that is,

$$u(k) = K_p T_D \frac{(e(k) - e(k-1)))}{T} \Rightarrow \boxed{D(z) = K_p T_D \frac{1 - z^{-1}}{T} = K_p T_D \frac{z - 1}{Tz}}$$

In many designs, the compensation is a sum of proportional and derivative control (or PD control). In this case, we have

$$D(z) = K_p \left(1 + \frac{T_D(z-1)}{Tz} \right).$$

or, equivalently,

$$\boxed{D(z) = K \frac{z - \alpha}{z}}$$

which is similar to the lead compensations that have been used in the designs in the previous sections. The difference is that the pole is at $z = 0$, whereas the pole has been placed at various locations along the z -plane real axis for the previous designs. In the continuous case, pure derivative control represents the ideal situation in that there is no destabilizing phase lag from the differentiation, or, equivalently, the pole is at $s = -\infty$. This s -plane pole maps into $z = 0$ for discrete rate control; however, the $z = 0$ pole does

add some phase lag because of the necessity to wait for one cycle in order to compute the first difference. Any other stable pole location, whether on the positive or negative real axis, would also have some delay or phase lag associated with it for the same reason.

5.8.3 Integral Control

For continuous systems, we integrate the error to arrive at the control,

$$u(t) = \frac{K_p}{T_I} \int_{t_0}^t e(t) dt \Rightarrow D(s) = \frac{K_p}{T_I s},$$

where T_I is called the *integral*, or *reset time*. The discrete equivalent is to sum all previous errors, yielding

$$u(k) = u(k-1) + \frac{K_p T}{T_I} e(k) \Rightarrow \boxed{D(z) = \frac{K_p T}{T_I(1-z^{-1})} = \frac{K_p T z}{T_I(z-1)}} \quad (5.60)$$

Just as for continuous systems, the primary reason for integral control is to reduce or eliminate steady-state errors, but this typically occurs at the cost of reduced stability.

5.8.4 PID Control

Combining all the above yields the PID controller

$$\boxed{D(z) = K_p \left(1 + \frac{Tz}{T_I(z-1)} + \frac{T_D(z-1)}{Tz} \right)}. \quad (5.61)$$

This form of control law is able satisfactorily to meet the specifications for a large portion of control problems and is therefore packaged commercially and sold for general use. The user simply has to determine the best values of K_p , T_D , and T_I .

5.8.5 Ziegler-Nichols PID Tuning

The parameters in the PID controller could be selected by any of the design methods previously discussed. However, these methods require a dynamic model of the process which is not always readily available. Ziegler-Nichols tuning is a method for picking the parameters based on fairly simple experiments on the process and thus bypasses the need to determine a complete dynamic model.

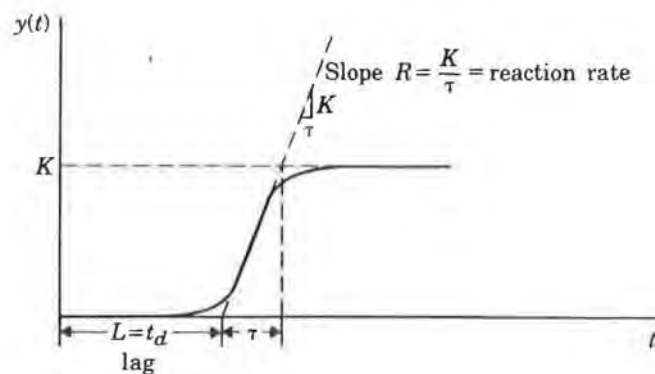


Figure 5.30 Process open-loop step response.

There are two methods. The first, called the *transient-response method*, requires that a step response of the open-loop system is obtained which looks something like that in Fig. 5.30. The response is reduced to two parameters, the time delay, L , and the steepest slope, R , which are defined in the figure. In order to achieve a damping of about $\zeta = 0.2$, the parameters are selected according to those in Table 5.2.

The second method is called the *stability-limit method*. The system is first controlled using proportional control only. The gain, K_p , is slowly increased until continuous oscillations result, at which point the gain and oscillation period are recorded and called K_u and P_u . The PID gains are then determined from Table 5.3.

The rules are based on continuous systems and will apply to the discrete case for very fast sampling (more than 20 times the bandwidth) provided the designer uses the value of T in (5.61) that reflects the actual sample period being used by the controller. For slower sampling, a response degradation similar to that in Example 5.3 should be expected, and additional rate control (higher T_D) would likely be required to make up for the sampling lag.

Table 5.2 Ziegler-Nichols tuning parameters using transient response.

	K_p	T_I	T_D
P	$1/RL$		
PI	$0.9/RL$	$3L$	
PID	$1.2/RL$	$2L$	$0.5L$

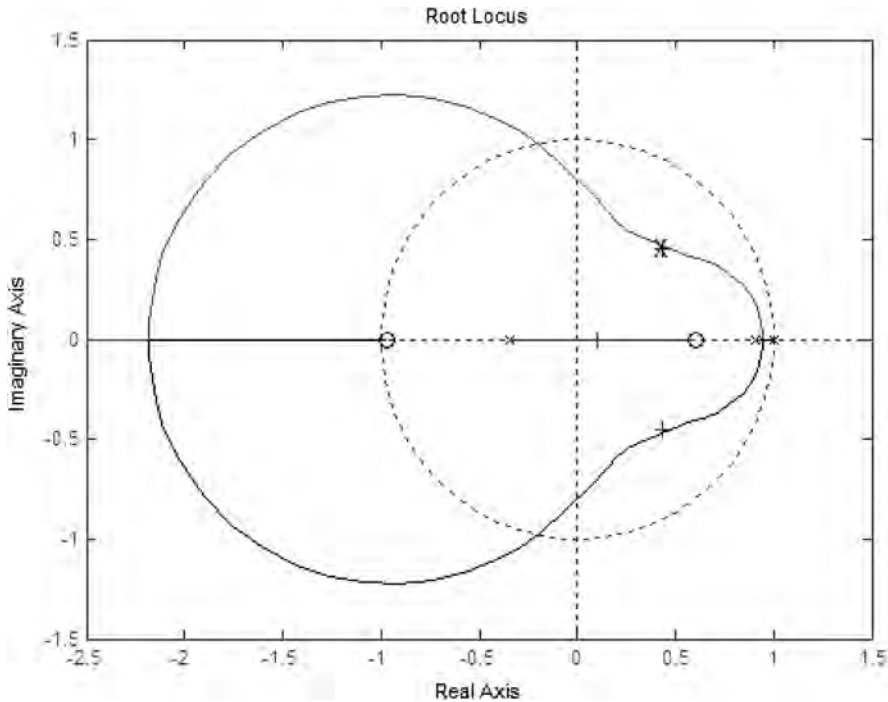


Figure 9.21 Root locus of the compensated system

Figure 9.21 shows the root locus of the compensated system. Clearly the locus passes through the required point. The d.c. gain at this point is $K = 123.9$.

The time response of the compensated system is shown in Figure 9.22.

9.2 PID CONTROLLER

The proportional–integral–derivative (PID) controller is often referred to as a ‘three-term’ controller. It is currently one of the most frequently used controllers in the process industry. In a PID controller the control variable is generated from a term proportional to the error, a term which is the integral of the error, and a term which is the derivative of the error.

- *Proportional*: the error is multiplied by a gain K_p . A very high gain may cause instability, and a very low gain may cause the system to drift away.
- *Integral*: the integral of the error is taken and multiplied by a gain K_i . The gain can be adjusted to drive the error to zero in the required time. A too high gain may cause oscillations and a too low gain may result in a sluggish response.
- *Derivative*: The derivative of the error is multiplied by a gain K_d . Again, if the gain is too high the system may oscillate and if the gain is too low the response may be sluggish.

Figure 9.23 shows the block diagram of the classical continuous-time PID controller. Tuning the controller involves adjusting the parameters K_p , K_d and K_i in order to obtain a satisfactory

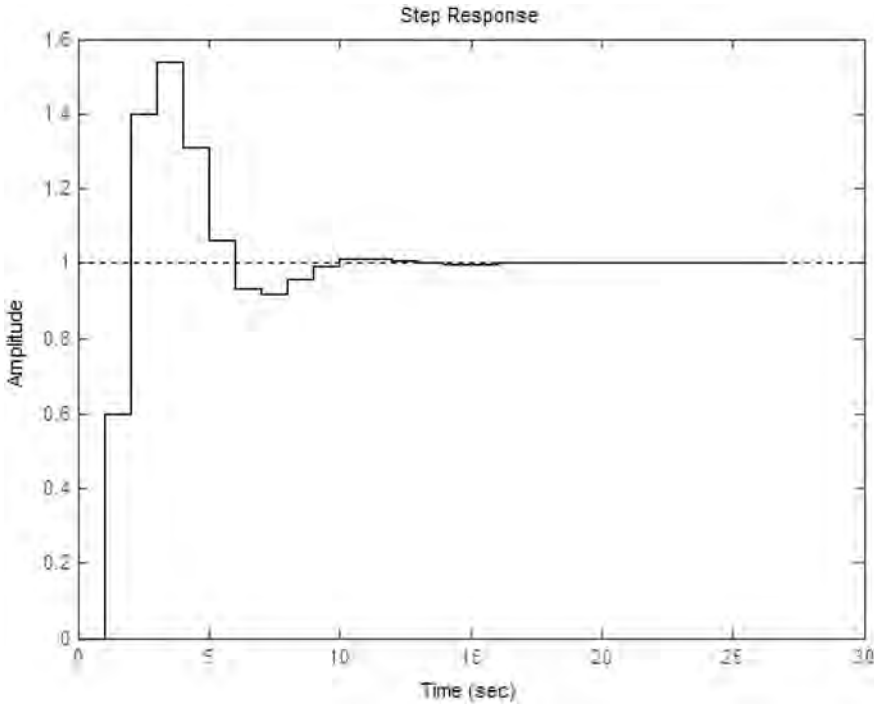


Figure 9.22 Time response of the system

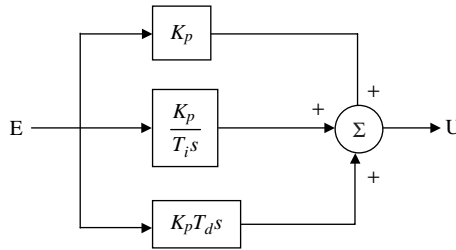


Figure 9.23 Continuous-time system PID controller

response. The characteristics of PID controllers are well known and well established, and most modern controllers are based on some form of PID.

The input–output relationship of a PID controller can be expressed as

$$u(t) = K_p \left[e(t) + \frac{1}{T_i} \int_0^t e(t) dt + T_d \frac{de(t)}{dt} \right], \quad (9.14)$$

where $u(t)$ is the output from the controller and $e(t) = r(t) - y(t)$, in which $r(t)$ is the desired set-point (reference input) and $y(t)$ is the plant output. T_i and T_d are known as the integral and

derivative action time, respectively. Notice that (9.14) is sometimes written as

$$u(t) = K_p e(t) + K_i \int_0^t e(t) dt + K_d \frac{de(t)}{dt} + u_0, \quad (9.15)$$

where

$$K_i = \frac{K_p}{T_i} \quad \text{and} \quad K_d = K_p T_d. \quad (9.16)$$

Taking the Laplace transform of (9.14), we can write the transfer function of a continuous-time PID as

$$\frac{U(s)}{E(s)} = K_p + \frac{K_p}{T_i s} + K_p T_d s. \quad (9.17)$$

To implement the PID controller using a digital computer we have to convert (9.14) from a continuous to a discrete representation. There are several methods for doing this and the simplest is to use the trapezoidal approximation for the integral and the backward difference approximation for the derivative:

$$\frac{de(t)}{dt} \approx \frac{e(kT) - e(kT - T)}{T} \quad \text{and} \quad \int_0^t e(t) dt \approx \sum_{k=1}^n T e(kT).$$

Equation (9.14) thus becomes

$$u(kT) = K_p \left[e(kT) + T_d \frac{e(kT) - e(kT - T)}{T} + \frac{T}{T_i} \sum_{k=1}^n e(kT) \right] + u_0. \quad (9.18)$$

The PID given by (9.18) is now in a suitable form which can be implemented on a digital computer. This form of the PID controller is also known as the *positional* PID controller. Notice that a new control action is implemented at every sample time.

The discrete form of the PID controller can also be derived by finding the z -transform of (9.17):

$$\frac{U(z)}{E(z)} = K_p \left[1 + \frac{T}{T_i(1 - z^{-1})} + T_d \frac{(1 - z^{-1})}{T} \right]. \quad (9.19)$$

Expanding (9.19) gives

$$\begin{aligned} u(kT) = & u(kT - T) + K_p [e(kT) - e(kT - T)] + \frac{K_p T}{T_i} e(kT) \\ & + \frac{K_p T_d}{T} [e(kT) - 2e(kT - T) + e(kT - 2T)]. \end{aligned} \quad (9.20)$$

This form of the PID controller is known as the *velocity* PID controller. Here the current control action uses the previous control value as a reference. Because only a change in the control action is used, this form of the PID controller provides a smoother bumpless control when the error is small. If a large error exists, the response of the velocity PID controller may be slow, especially if the integral action time T_i is large.

The two forms of the PID algorithm, (9.18) and (9.20), may look quite different, but they are in fact similar to each other. Consider the positional controller (9.18). Shifting back one

sampling interval, we obtain

$$u(kT - T) = K_p \left[e(kT - T) + T_d \frac{e(kT - T) - e(kT - 2T)}{T} + \frac{T}{T_i} \sum_{k=1}^{n-1} e(kT) \right] + u_0.$$

Subtracting from (9.18), we obtain the velocity form of the controller, as given by (9.20).

9.2.1 Saturation and Integral Wind-Up

In practical applications the output value of a control action is limited by physical constraints. For example, the maximum voltage output from a device is limited. Similarly, the maximum flow rate that a pump can supply is limited by the physical capacity of the pump. As a result of this physical limitation, the error signal does not return to zero and the integral term keeps adding up continuously. This effect is called integral wind-up (or integral saturation), and as a result of it long periods of overshoot can occur in the plant response. A simple example of what happens is the following. Suppose we wish to control the position of a motor and a large set-point change occurs, resulting in a large error signal. The controller will then try to reduce the error between the set-point and the output. The integral term will grow by summing the error signals at each sample and a large control action will be applied to the motor. But because of the physical limitation of the motor electronics the motor will not be able to respond linearly to the applied control signal. If the set-point now changes in the other direction, then the integral term is still large and will not respond immediately to the set-point request. Consequently, the system will have a poor response when it comes out of this condition.

The integral wind-up problem affects positional PID controllers. With velocity PID controllers, the error signals are not summed up and as a result integral wind-up will not occur, even though the control signal is physically constrained.

Many techniques have been developed to eliminate integral wind-up from the PID controllers, and some of the popular ones are as follows:

- Stop the integral summation when saturation occurs. This is also called conditional integration. The idea is to set the integrator input to zero if the controller output is saturated and the input and output are of the same sign.
- Fix the limits of the integral term between a minimum and a maximum.
- Reduce the integrator input by some constant if the controller output is saturated. Usually the integral value is decreased by an amount proportional to the difference between the unsaturated and saturated (i.e maximum) controller output.
- Use the velocity form of the PID controller.

9.2.2 Derivative Kick

Another possible problem when using PID controllers is caused by the derivative action of the controller. This may happen when the set-point changes sharply, causing the error signal to change suddenly. Under such a condition, the derivative term can give the output a *kick*, known as a *derivative kick*. This is usually avoided in practice by moving the derivative term

to the feedback loop. The proportional term may also cause a sudden kick in the output and it is common to move the proportional term to the feedback loop.

9.2.3 PID Tuning

When a PID controller is used in a system it is important to tune the controller to give the required response. Tuning a PID controller involves selecting values for the controller parameters K_p , T_i and T_d . There are many techniques for tuning a controller, ranging from the first techniques described by J.G. Ziegler and N.B. Nichols (known as the Ziegler–Nichols tuning algorithm) in 1942 and 1943, to recent auto-tuning controllers. In this section we shall look at the tuning of PID controllers using the Ziegler–Nichols tuning algorithm.

Ziegler and Nichols suggested values for the PID parameters of a plant based on open-loop or closed-loop tests of the plant. According to Ziegler and Nichols, the open-loop transfer function of a system can be approximated with a time delay and a single-order system, i.e.

$$G(s) = \frac{K e^{-sT_D}}{1 + sT_1}, \quad (9.21)$$

where T_D is the system time delay (i.e. transportation delay), and T_1 is the time constant of the system.

9.2.3.1 Open-Loop Tuning

For open-loop tuning, we first find the plant parameters by applying a step input to the open-loop system. The plant parameters K , T_D and T_1 are then found from the result of the step test as shown in Figure 9.24.

Ziegler and Nichols then suggest using the PID controller settings given in Table 9.1 when the loop is closed. These parameters are based on the concept of minimizing the integral of the absolute error after applying a step change to the set-point.

An example is given below to illustrate the method used.

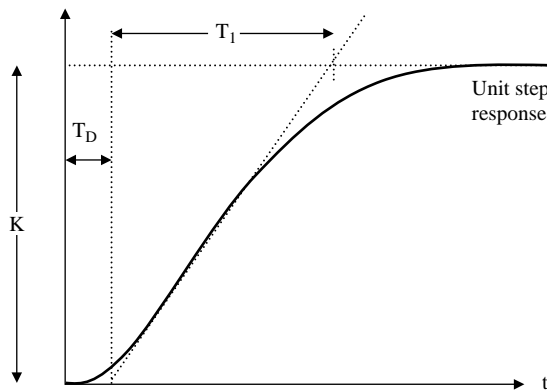


Figure 9.24 Finding plant parameters K , T_D and T_1

Table 9.1 Open-loop Ziegler–Nichols settings

Controller	K_p	T_i	T_d
Proportional	$\frac{T_1}{KT_D}$		
PI	$\frac{0.9T_1}{KT_D}$	$3.3T_D$	
PID	$\frac{1.2T_1}{KT_D}$	$2T_D$	$0.5T_D$

Example 9.7

The open-loop unit step response of a thermal system is shown in Figure 9.25. Obtain the transfer function of this system and use the Ziegler–Nichols tuning algorithm to design (a) a proportional controller, (b) to design a proportional plus integral (PI) controller, and (c) to design a PID controller. Draw the block diagram of the system in each case.

Solution

From Figure 9.25, the system parameters are obtained as $K = 40^\circ\text{C}$, $T_D = 5$ s and $T_1 = 20$ s, and the transfer function of the plant is

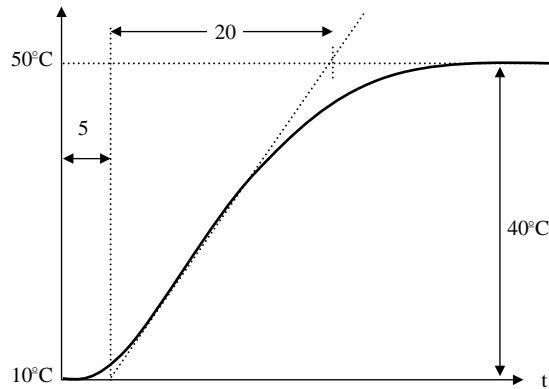
$$G(s) = \frac{40e^{-5s}}{1 + 20s}.$$

Proportional controller. According to Table 9.1, the Ziegler–Nichols settings for a proportional controller are:

$$K_p = \frac{T_1}{KT_D}.$$

Thus,

$$K_p = \frac{20}{40 \times 5} = 0.1,$$

**Figure 9.25** Unit step response of the system for Example 9.7

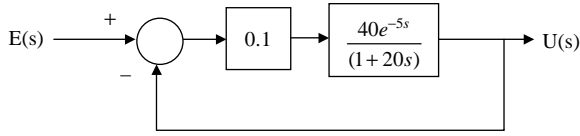


Figure 9.26 Block diagram of the system with proportional controller

The transfer function of the controller is then

$$\frac{U(s)}{E(s)} = 0.1,$$

and the block diagram of the closed-loop system with the controller is shown in Figure 9.26.

PI controller. According to Table 9.1, the Ziegler–Nichols settings for a PI controller are

$$K_p = \frac{0.9T_1}{KT_D} \quad \text{and} \quad T_i = 3.3T_D.$$

Thus,

$$K_p = \frac{0.9 \times 20}{40 \times 5} = 0.09 \quad \text{and} \quad T_i = 3.3 \times 5 = 16.5.$$

The transfer function of the controller is then

$$\frac{U(s)}{E(s)} = 0.09 \left[1 + \frac{1}{16.5s} \right] = \frac{0.09(16.5s + 1)}{16.5s}$$

and the block diagram of the closed-loop system with the controller is shown in Figure 9.27.

PID controller. According to Table 9.1, the Ziegler–Nichols settings for a PID controller are

$$K_p = \frac{1.2T_1}{KT_D}, \quad T_i = 2T_D, \quad T_d = 0.5T_D.$$

Thus,

$$K_p = \frac{1.2 \times 20}{40 \times 5} = 0.12, \quad T_i = 2 \times 5 = 10, \quad T_d = 0.5 \times 5 = 2.5.$$

The transfer function of the required PID controller is

$$\frac{U(s)}{E(s)} = K_p \left[1 + \frac{1}{T_i s} + T_d s \right] = 0.12 \left[1 + \frac{1}{10s} + 2.5s \right]$$

or

$$\frac{U(s)}{E(s)} = \frac{3s^2 + 1.2s + 0.12}{10s}.$$

The block diagram of the system, together with the controller, is shown in Figure 9.28.

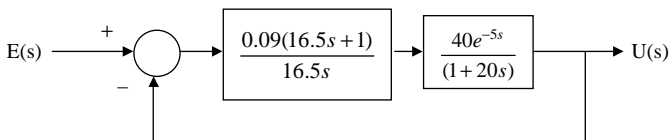


Figure 9.27 Block diagram of the system with PI controller

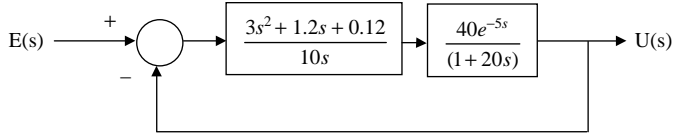


Figure 9.28 Block diagram of the system with PID controller

9.2.3.2 Closed-Loop Tuning

The Ziegler–Nichols closed-loop tuning algorithm is based on plant closed-loop tests. The procedure is as follows:

- Disable any derivative and integral action in the controller and leave only the proportional action.
- Carry out a set-point step test and observe the system response.
- Repeat the set-point test with increased (or decreased) controller gain until a stable oscillation is achieved (see Figure 9.29). This gain is called the *ultimate gain*, K_u .
- Read the period of the steady oscillation and let this be P_u .
- Calculate the controller parameters according to the following formulae: $K_p = 0.45K_u$, $T_i = P_u/1.2$ in the case of the PI controller; and $K_p = 0.6K_u$, $T_i = P_u/2$, $T_d = T_u/8$ in the case of the PID controller.

9.3 EXERCISES

1. The open-loop transfer function of a plant is given by:

$$G(s) = \frac{e^{-4s}}{1 + 2s}.$$

- (a) Design a dead-beat digital controller for the system. Assume that $T = 1$ s.
- (b) Draw the block diagram of the system together with the controller.
- (c) Plot the time response of the system.

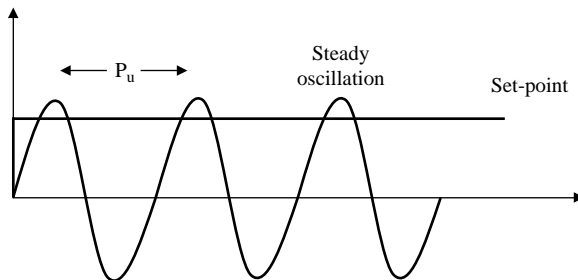


Figure 9.29 Ziegler–Nichols closed-loop test