# Effects of Quantization

- **Floating point vs. fixed point**

- **Analysis of effects of quantization on *parameters***

- **Analysis of effects of quantization on *signals***

  - **Worst case**

  - **Steady-state worst case**

  - **Example**

# Effects of Quantization

**When a digital computer is used for controlling a plant, the control $u(k)$ can:**

    **1. Only change at discrete times**
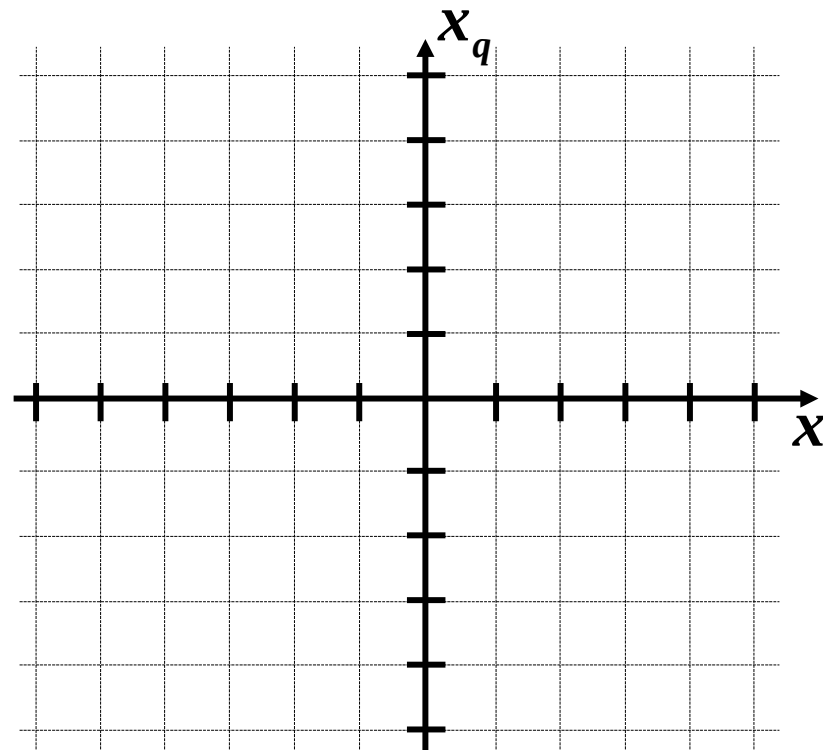    **2. Only be represented to finite accuracy**

## Quantization:

  – **Required by digital microprocessor**

$$x \longrightarrow \boxed{\textbf{A/D}} \longrightarrow x_q$$

$$q = 2^{-\ell}$$

**where** $\ell$ **is the number of bits**

# Floating Point vs. Fixed Point

## Floating point

**Advantage:  more flexible, more accurate**

**Disadvantage:  slower, more difficult, more costly**

## Fixed point

**Advantage:  economical, faster, easier**

**Disadvantage:  less flexible, underflow and overflow,
                               less  accurate, scaling needed.**

# Recommended Implementation Architectures

**Brief discussion of quantization of *parameters* :**

$$u \longrightarrow \boxed{\frac{b_1 z^2 + b_2 z + b_3}{z^2 + a_1 z + a_2}} \longrightarrow y$$

$$u \longrightarrow \boxed{\frac{(b_1 + \varepsilon_{b1})z^2 + (b_2 + \varepsilon_{b2})z + (b_3 + \varepsilon_{b3})}{z^2 + (a_1 + \varepsilon_{a1})z + (a_2 + \varepsilon_{a2})}} \longrightarrow y$$

# Effect of Parameter Storage Errors

**In general, characteristic equation:**

$$P(z, \ \alpha) = z^n + \alpha_1 z^{n-1} + \cdots + \alpha_k z^{n-k} + \cdots + \alpha_n$$
$$= (z - \lambda_1)(z - \lambda_2) \cdots (z - \lambda_j) \cdots (z - \lambda_n) = 0$$

**When there are no word length limitations in storing the coefficients** $\alpha_k$**, the nominal roots are** $z = \lambda_j$

**Suppose there is an error in a coefficient:** $\quad \alpha_k \rightarrow \alpha_k + \delta\alpha_k$

**This causes an error in the poles:** $\quad \lambda_j \rightarrow \lambda_j + \delta\lambda_j$

**Stability can be affected!**

$$P\left(\lambda_j + \delta\lambda_j, \quad \alpha_k + \delta\alpha_k\right) = 0$$

**Calculate sensitivity:** $\dfrac{\delta\lambda_j}{\delta\alpha_k}$

$$P\left(\lambda_j + \delta\lambda_j, \quad \alpha_k + \delta\alpha_k\right) = P\left(\lambda_j, \quad \alpha_k\right) + \left.\frac{\partial P}{\partial z}\right|_{\substack{z=\lambda_j\\\alpha}} \delta\lambda_j + \left.\frac{\partial P}{\partial\alpha_k}\right|_{\substack{z=\lambda_j\\\alpha}} \delta\alpha_k$$

$$+ \quad \textbf{higher} \quad \textbf{order} \quad \textbf{terms} = 0$$

$$\left.\frac{\partial P}{\partial z}\right|_{\substack{z=\lambda_j\\\alpha}} \delta\lambda_j + \left.\frac{\partial P}{\partial\alpha_k}\right|_{\substack{z=\lambda_j\\\alpha}} \delta\alpha_k = 0$$

**Sensitivity:** $\dfrac{\partial\lambda_j}{\partial\alpha_k} = \left.\dfrac{-\,\partial P / \partial\alpha_k}{\partial P / \partial z}\right|_{\substack{z=\lambda_j\\\alpha}}$
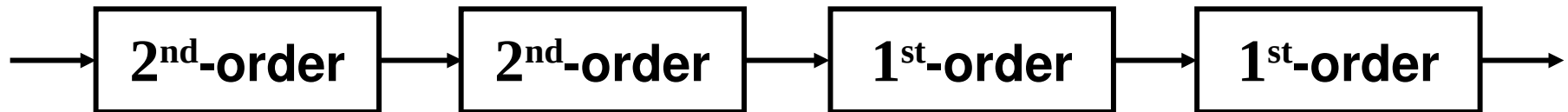
# Sensitivity for direct realization:

$$\delta\lambda_j = \frac{-\lambda_j^{n-k}}{\prod_{\ell \neq j}(\lambda_j - \lambda_\ell)}\delta\alpha_k$$

1. **All poles are assumed to be distinct.  If any roots have multiplicities greater than $1$, then the sensitivities of these roots are infinite.**

2. **Pole locations are especially sensitive to parameter changes/errors if poles are close.**

3. **In direct realization, the sensitivity to parameter storage errors are higher than <u>parallel</u> or <u>cascade</u> realization.**
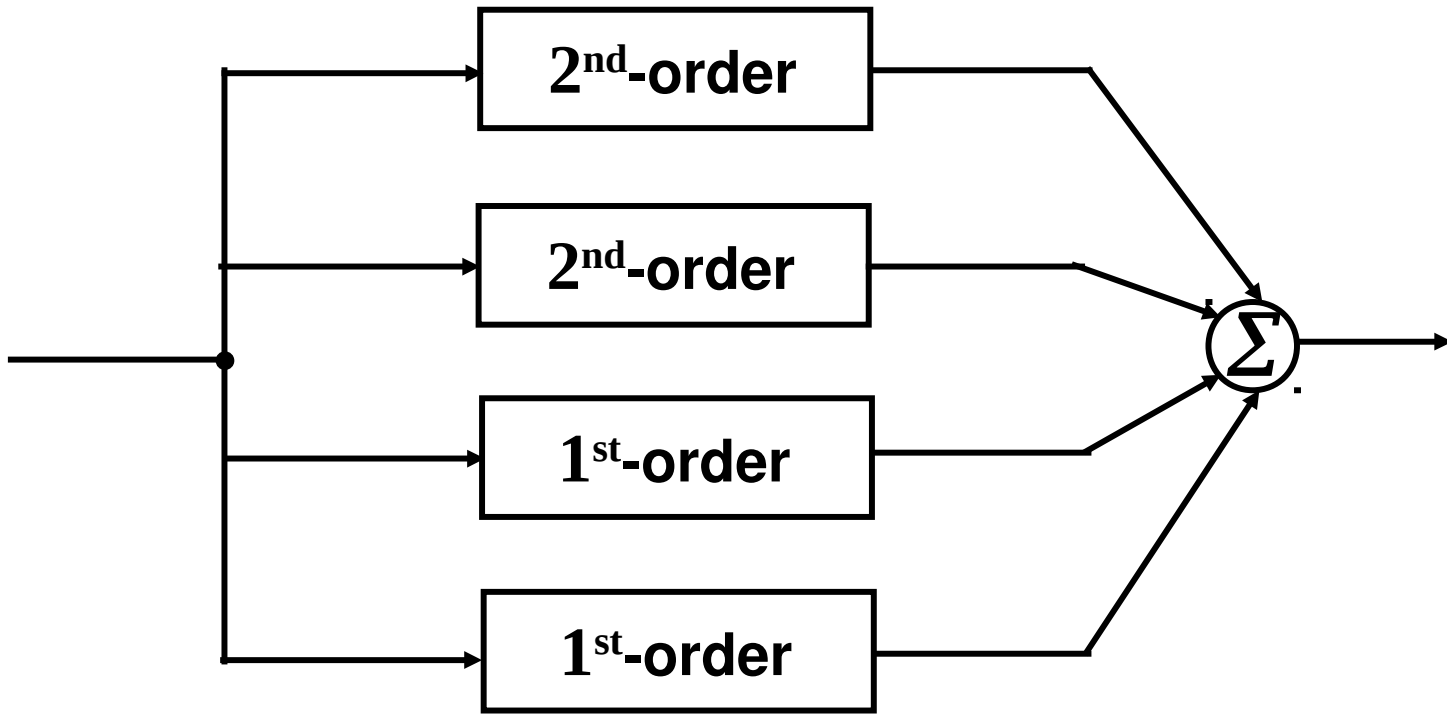
**This brief analysis indicates that the best way (i.e., least sensitive to quantization of parameters) to implement a higher-order controller is to break it down into first and second-order components and use either a *cascade* or *parallel* implementation.**
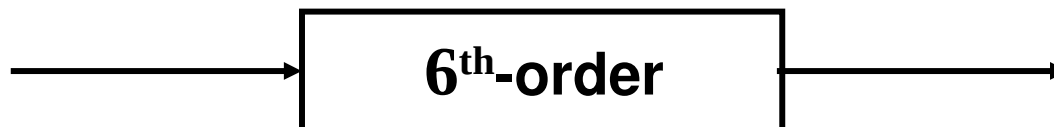
**Example:**

A $6^{th}$-order system with $2$ pairs of complex poles and $2$ real poles can be implemented with a **cascade** of lower-order factors:

| $2^{nd}$-order | $2^{nd}$-order | $1^{st}$-order | $1^{st}$-order |

# Or, in a parallel implementation:

```
                    ┌──────────────┐
         ┌─────────▶│   2nd-order  │──────────┐
         │          └──────────────┘          │
         │          ┌──────────────┐          │
         │  ┌──────▶│   2nd-order  │────────┐ │
         │  │       └──────────────┘        ▼ ▼
──────▶──●──┤                              ╭─────╮
            │       ┌──────────────┐       │  Σ  │──────▶
            ├──────▶│   1st-order  │──────▶│     │
            │       └──────────────┘       ╰─────╯
            │       ┌──────────────┐        ▲
            └──────▶│   1st-order  │────────┘
                    └──────────────┘
```

**Using a direct realization or a direct implementation of the $6^{th}$ order controller can be significantly more sensitive to quantization of parameters:**

```
              ┌──────────────┐
─────────────▶│   6th-order  │──────────────▶
              └──────────────┘
```

# Now let's analyze the effects of quantization on *signals* :

$$u \longrightarrow \boxed{\phantom{xxxxxxxxxxxx}} \longrightarrow y$$
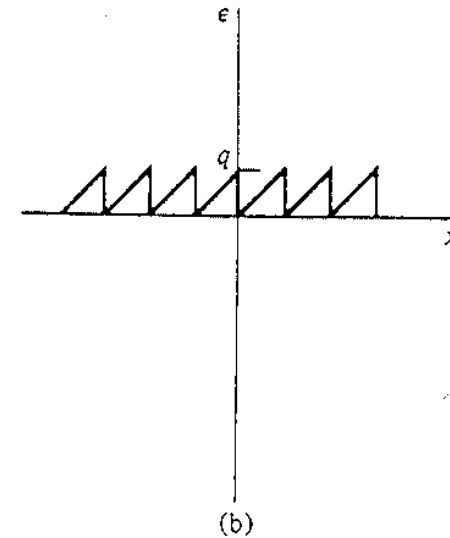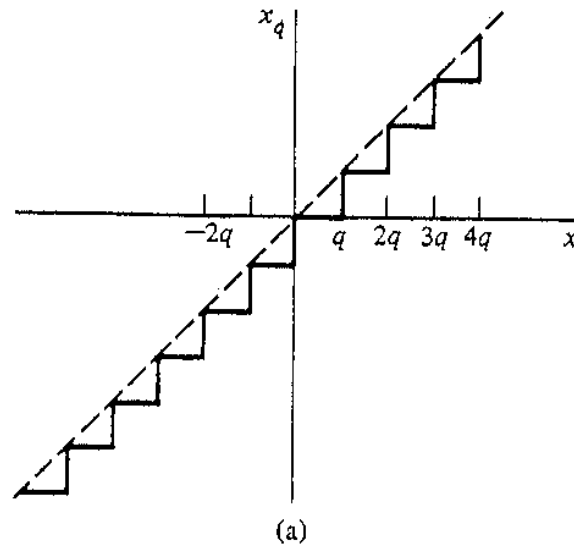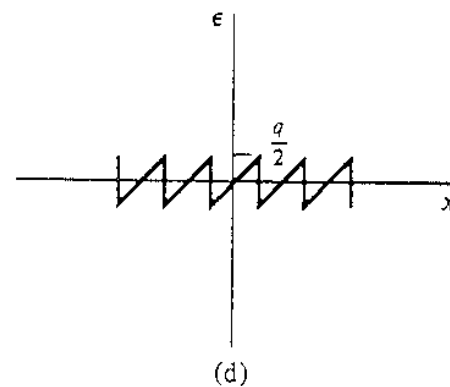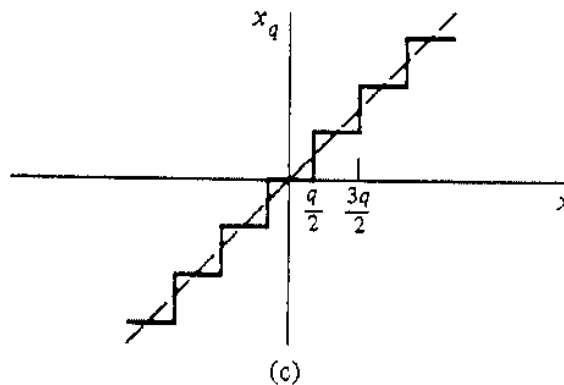
# Effect of Quantization on Signals

**Figure 10.1**
Plot of effects of number truncation. (a) Plot of variable versus truncated values. (b) Plot of error due to truncation. (c) Plot of variable versus rounded values. (d) Round-off error
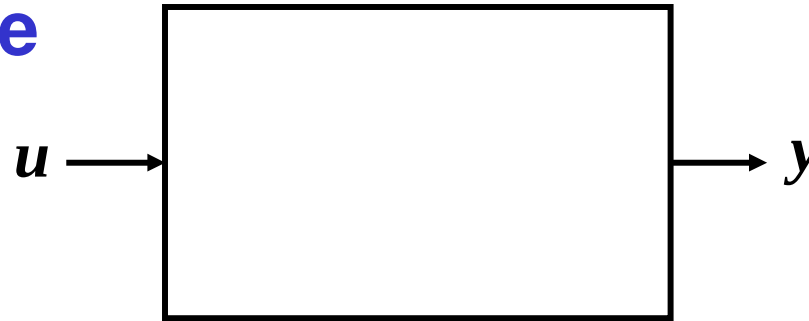
**Truncation:**

**Rounding:**

# How Do We Represent $\varepsilon$ ?

**Difficult to do exactly because of the nonlinear nature of $\varepsilon$ .**

**Assume fixed-point and round-off.**

1. **Worst case** (Bertram): $\varepsilon$ is selected to make effect of error as large as possible.

2. **Steady-state worst case** (Slaughter and Blackman): for round-off,

$$\varepsilon \equiv \frac{q}{2} = \text{constant}$$

# 1.  Worst Case

$$u \longrightarrow \boxed{\phantom{xxxxxxxxxxxx}} \longrightarrow y$$

**Ideally:**     $\xi(k+1) = \Phi\xi(k) + \Gamma u(k)$

$\qquad\qquad y(k) = H\xi(k)$

**In reality:**     $\hat{\xi}(k+1) = \Phi\hat{\xi}(k) + \Gamma u(k) + \Gamma_1 \varepsilon(k, x)$

$\qquad\qquad \hat{y}(k) = H\hat{\xi}(k)$

**Let:**     $\tilde{\xi}(k) = \xi(k) - \hat{\xi}(k), \qquad \tilde{y}(k) = y(k) - \hat{y}(k)$

**Then:**     $\tilde{\xi}(k+1) = \Phi\tilde{\xi}(k) - \Gamma_1 \varepsilon(k, x)$

$\qquad\qquad \tilde{y}(k) = H\tilde{\xi}(k)$

**For round-off:** $\quad |\varepsilon| \le \dfrac{q}{2}$

**The transfer function from** $\quad \varepsilon \to \tilde{y}:$

$$\tilde{Y}(z) = -\mathbf{H}(z\mathbf{I} - \mathbf{\Phi})^{-1}\mathbf{\Gamma}_1 E(z, x)$$

**Inverse $Z$-transform:** $\quad \tilde{y}(k) = \sum_{n=0}^{k} h_1(n)\varepsilon(k - n)$

**Assuming $\varepsilon$ is selected so that** $\quad \tilde{y}(k)$ **is as large as possible.**

$$\left|\tilde{y}(k)\right| = \left|\sum_{n=0}^{k} h_1(n)\varepsilon(k - n)\right|$$

$$\boxed{\left|\tilde{y}(k)\right| \le \dfrac{q}{2}\sum_{n=0}^{\infty} \left|h_1\right|}$$ **Worst case bound (pessimistic)**

# 2.  Steady-State Worst Case

$$\tilde{y}(k) = \sum_{n=0}^{k} h_1(n)\varepsilon(k - n)$$

**Assume  $\varepsilon \equiv \dfrac{q}{2}$  for all the time.**

$$\left|\tilde{y}_{ss}\right| = \left|\tilde{y}(\infty)\right| = \frac{q}{2}\left|\sum_{n=0}^{\infty} h_1(n)\right|$$

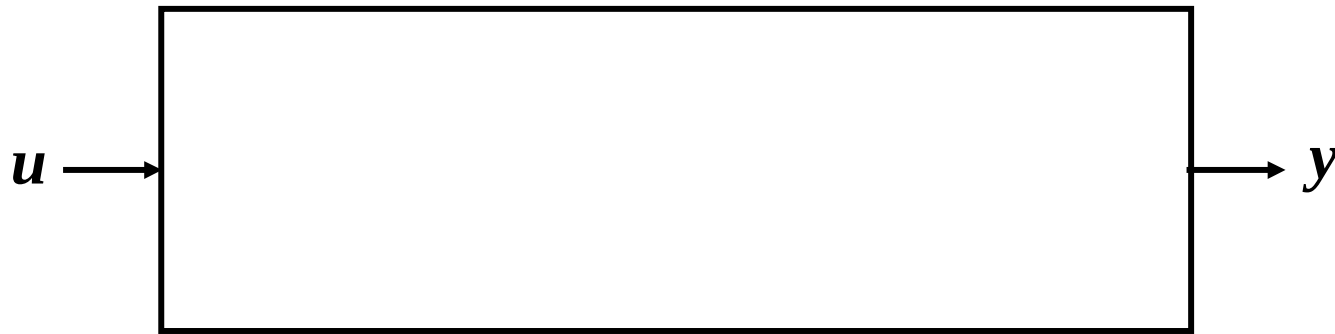$$H_1(z) = \sum_{n=0}^{\infty} h_1(n)z^{-n}$$

# Comparison

**Bertram's Worst Case Bound:**

$$|\tilde{y}| \leq \sum |h| \frac{q}{2}$$

**Steady-State Worst Case:**

$$|\tilde{y}_{ss}| \approx \left|\sum h\right| \frac{q}{2}$$

# Generalization to Multiple Sources of Round-off Error

$$u \longrightarrow \boxed{\phantom{XXXXXXXXXXXXXXX}} \longrightarrow y$$

**Bertram's worst case:**

$$|\tilde{y}| \leq \frac{q_1}{2}\sum_{n=0}^{\infty}|h_1(n)| + \frac{q_2}{2}\sum_{n=0}^{\infty}|h_2(n)| + \cdots$$

$$|\tilde{y}| \leq \sum_{j=1}^{K}\frac{q_j}{2}\sum_{n=0}^{\infty}|h_j(n)|$$

**If all $q_j$'s are equal:**

$$|\tilde{y}| \leq \frac{q}{2}\sum_{j=1}^{K}\sum_{n=0}^{\infty}|h_j(n)|$$

**Steady-state worst case:**

$$|\tilde{y}_{ss}| \approx \frac{q_1}{2}|H_1(1)| + \frac{q_2}{2}|H_2(1)| + \cdots$$

$$|\tilde{y}_{ss}| \approx \sum_{j=1}^{K}\frac{q_j}{2}|H_j(1)|$$

**If all $q_j$'s are equal:**

$$|\tilde{y}_{ss}| \approx \frac{q}{2}\sum_{j=1}^{K}|H_j(1)|$$

# Quantization Example

## Example 10.4 of Text

**Second-order system:**

$$y(k + 2) + a_1 y(k + 1) + a_2 y(k) = \frac{1 + a_1 + a_2}{2}\left[u(k + 1) + u(k)\right]$$

**Compute error at $y$ using:**     1)  **Worst-case bound**
                                     2)  **Steady-state estimate**

2)  **Steady-state estimate**

$$\left|\tilde{y}_{ss}\right| \approx \frac{q}{2}\left|H\,(1)\right|$$

Assume quantization error $\varepsilon$ enters at same
point where input $u$ enters system.          $H\,(z) = \ ?$

$$H(z) = \frac{Y(z)}{U(z)} = \frac{1 + a_1 + a_2}{2} \frac{z + 1}{z^2 + a_1 z + a_2}$$

$$\left| \tilde{y}_{ss} \right| \approx \frac{q}{2}$$

**Normalized error:**          $\left| \dfrac{\tilde{y}_{ss}}{q/2} \right| \approx 1$

# 1)  Worst-case bound

$$\left|\tilde{y}_{ss}\right| \leq \frac{q}{2} \sum \left|h\right|$$

**Normalized:**

$$\left|\frac{\tilde{y}_{ss}}{q/2}\right| \leq \sum_{n=0}^{\infty} \left|h\right|$$

**What is $h$ ?**

$$h = Z^{-1}\left\{H(z)\right\} = Z^{-1}\left[\frac{1 + a_1 + a_2}{2} \frac{z + 1}{z^2 + a_1 z + a_2}\right]$$
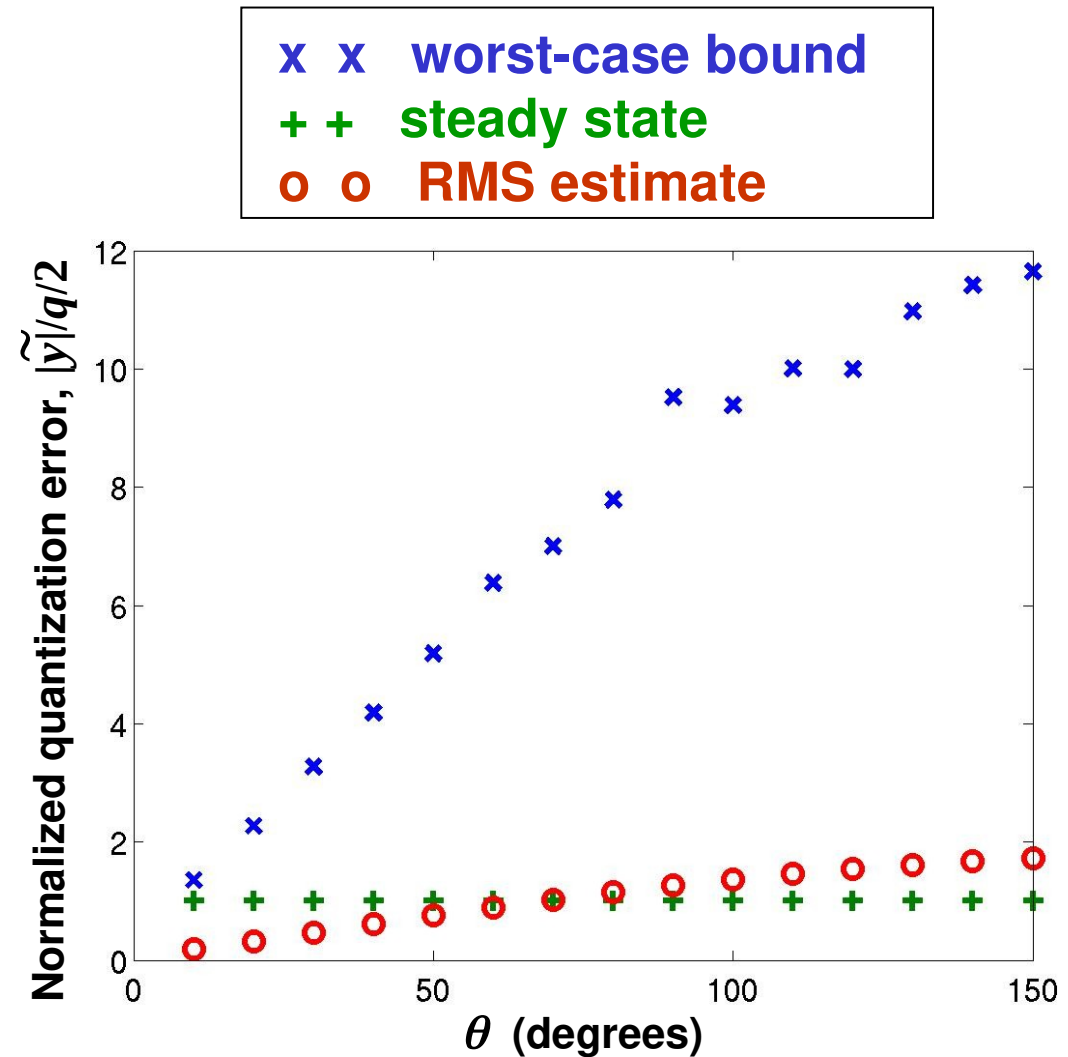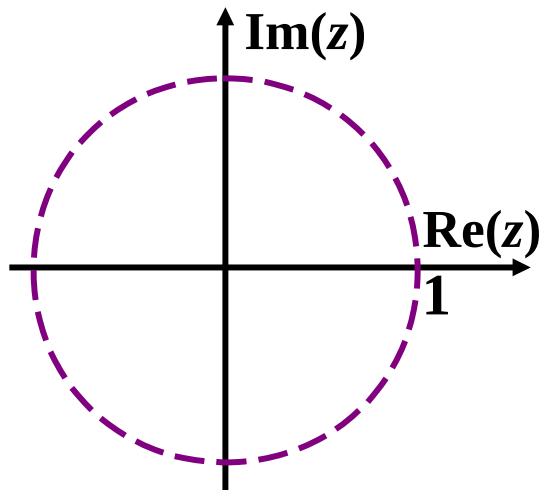
**For**   $a_1 = -2r\cos\theta,$   $a_2 = r^2$   $\implies$   **poles at**   $re^{\pm j\theta}$

**Using $Z$-transform table:**

It is difficult to compute $\displaystyle\sum_{k=0}^{\infty}\left|h(k)\right|$ analytically, but it is easy to write a script in Matlab to compute $\displaystyle\sum_{k=0}^{N}\left|h(k)\right|$ where $N$ is such that $\left|h(k)\right|$ is "small".

**poles at** $re^{\pm j\theta}$

$r = 0.9$

**Im(z)**

**Re(z)**

**1**

**x  x   worst-case bound**

**+  +   steady state**

**o  o   RMS estimate**



Normalized quantization error, $|\tilde{y}|/q/2$

$\theta$ **(degrees)**

# Other Quantization Topics

- **Stochastic analysis of quantization error (Section 10.1)**


- **Limit cycles and dither (Section 10.3)**