

25 YEARS ANNIVERSARY

SOKT

The central graphic element is a large, white, stylized number "25" with a circular arc above it containing the text "YEARS ANNIVERSARY". Below the "25", the letters "SOKT" are written in a bold, white, sans-serif font.

**ĐẠI HỌC BÁCH KHOA HÀ NỘI
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG**



ĐẠI HỌC BÁCH KHOA HÀ NỘI
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

APPLIED STATISTICS AND EXPERIMENTAL DESIGN INTRODUCTION

Applied Statistics and Experimental Design

Assoc. Prof. Nguyễn Linh Giang
School of Information and Communication
Technology
giangnl@soict.hust.edu.vn

Applied Statistics and Experimental Design

- Course description
 - For undergraduated and graduated student.
 - Build models of stochastic processes
 - Analysis of uncertainty
 - Statistical inference
 - Design of Experiments and Analysis of experimental data

Applied Statistics and Experimental Design

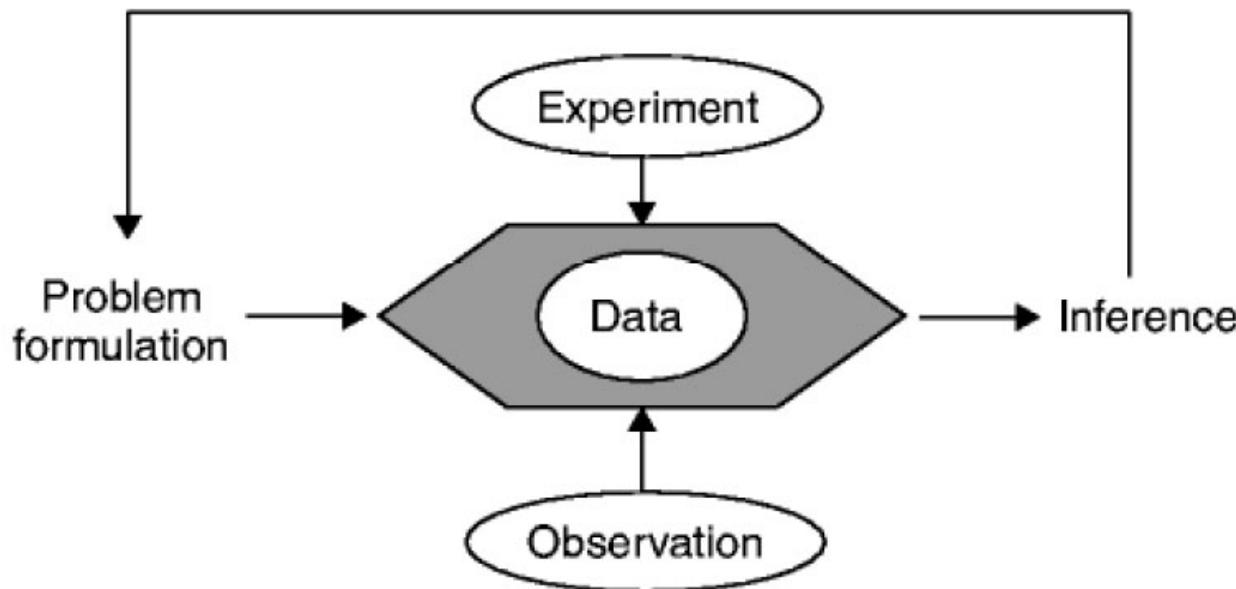
- Grading
 - Project/experimental exercise/Midterm test: 30%
 - Final exam: 70%
- Requisite
 - Mathematical Analysis, Calculus.
 - Probability theory and statistics.

Course outline

- Introduction to Applied Statistics and Data Analysis with Experiments
- Models and Definitions
- Basics of Probability Theory
- Basics of Statistics
- Random Processes
 - Time series analysis
- Statistical Errors and Estimation
- Data Analysis and Experimental Design

Applied Statistics and Experimental Design

- Stages in statistical analysis of random data



References

- J. S. Bendat, A. G. Piersol. Random Data: analysis and measurement procedures.
- Trossets M. W, An introductions to statistical inference and data analysis.
- Papoulis, Probability, random variables and stochastic processes.
- Walpole, Probability and Statistics for Engineers and Scientists

Statistics and Data Analysis

- Variability in Scientific Data
 - The use of statistical methods in many areas involves the gathering of information or scientific data.
 - Data have been collected, summarized, reported, and stored for perusal
 - Distinction between collection of scientific information and inferential statistics.
- Recent attention – inferential statistics => statistical methods employed by statistical practitioners
 - Statistical methods: scientific judgment in the face of **uncertainty** and **variation**.

Statistics and Data Analysis

- Statistical methods are used to analyze data from a process to improve the **quality** of the process.
- **Inferential statistics** - using analytical methods that allow us to go beyond merely reporting data to drawing conclusions (or inferences) about the scientific system.
- Statisticians make use of fundamental laws of probability and statistical inference to draw conclusions about scientific systems.

Statistics and Data Analysis

- Information is gathered in the form of **samples**, or collections of **observations**.
- Samples are collected from **populations** - collections of all individuals or individual items of a particular type.
- At times a population signifies a scientific system.
 - For example: to eliminate defects in a manufacturing computer boards .
 - A sampling process: collecting information on 50 computer boards sampled randomly from the process.
 - The population - all computer boards manufactured
 - After an improvement in the computer board process
 - a second sample of boards is collected,
 - any conclusions drawn regarding the effectiveness of the change in process should extend to the entire population of computer boards produced under the “improved process.”

Statistics and Data Analysis

- Methods of Statistical Inference
 - Experimental Study: Collect scientific data in a systematic way: with planning.
 - At times the planning is, by necessity, quite limited.
 - The factors - certain properties or characteristics of the items or objects in the population.
 - Experimental Design: move the factors to different levels according to prescription
 - Observational study: data are collected in the field but factor levels can not be preselected.
- In the former, the quality of the inferences will depend on proper planning of the experiment.
- In the latter, the scientist is at the mercy of what can be gathered.

Statistics and Data Analysis

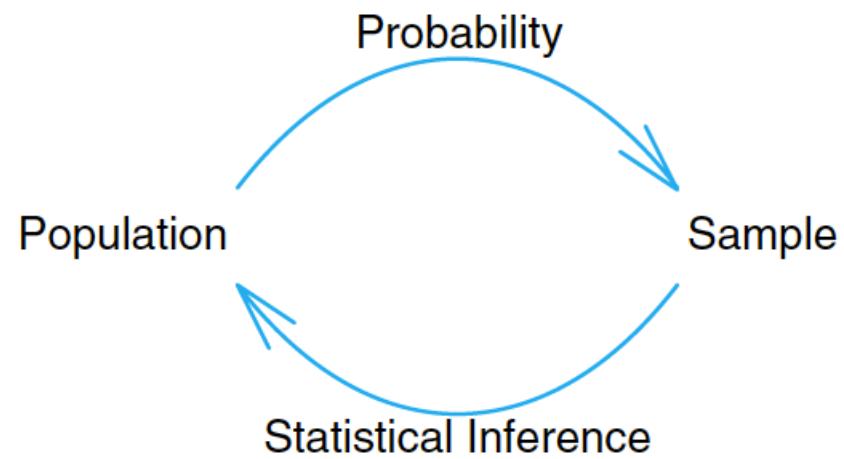
- The Importance of statistical thinking:
 - Research scientists gain much from scientific data
 - Data provide understanding of scientific phenomena
 - Product and process engineers learn a great deal in their off-line efforts to improve the process

Statistics and Data Analysis

- **Descriptive statistics:** a set of single-number statistics or descriptive statistics – when inferential statistics is not required - some sort of summary of a set of data represented in the sample.
 - These numbers give a sense of center of the location of the data, variability in the data, and the general nature of the distribution of observations in the sample.
 - Descriptive statistics are accompanied by graphics, tools for computation of important characteristics of the data set (mean, medians, standard deviation, ...)

Statistica and Data Analysis

- Probability and Inferential Statistics
 - Inductive reasoning
 - Deductive reasoning
- The bridge between the data and the conclusion: is based on foundations of statistical inference, distribution theory, and sampling distributions



Statistics and Data Analysis

Sampling procedure and Collection of Data

- Simple Random Sampling
 - Assumption: only a single population exists in the problem.
 - Simple random sampling – any sample of a specified sample size are equiprobable to be selected as any other sample of the same size.
 - Sample size - the number of elements in the sample.
 - The biased samples – the samples chosen in some limited population

Statistics and Data Analysis

Sampling procedure and Collection of Data

- Stratified Random Sampling

- The sampling units are not homogeneous and naturally divide themselves into nonoverlapping homogeneous groups called *strata* or *classes*.
- Stratified random sampling - random selection of a sample within each stratum.

Statistics and Data Analysis

Sampling procedure and Collection of Data

- Experimental Design

- The concept of randomness or random assignment plays a huge role in the area of experimental design
- A set of so-called treatments or treatment combinations becomes the populations to be studied or compared in some sense.
- Variability in the experimental unit – important concept in inferential statistics
- Completely randomized design
- Measure of variability – descriptive statistics

Statistics and Data Analysis

Measures of Locations

- Sample means and Median
 - Quantitative values of where the center of data is located
 - Suppose that the observations in a sample are x_1, x_2, \dots, x_n
 - Sample mean: denoted by \bar{x} - centroid of the data in a sample
$$\bar{x} = \sum_{i=1}^n \frac{x_i}{n} = \frac{1}{n}(x_1 + x_2 + \dots + x_n)$$
 - The basis for the computation of \bar{x} is that of an estimate of the population mean
 - Sample median: arrange the observations in increasing order: x_1, x_2, \dots, x_n , denoted by \tilde{x}
 - $\tilde{x} = \begin{cases} x_{(n+1)/2} & \text{if } n \text{ is odd} \\ \frac{1}{2}(x_{\frac{n}{2}} + x_{\frac{n}{2}+1}) & \text{if } n \text{ is even} \end{cases}$

Statistics and Data Analysis

Measures of Locations

- The other measures of location
 - There are several other methods of quantifying the center of location of the data in the sample.
 - Alternatives to the sample mean are designed to produce values that represent compromises between the mean and the median
 - Trimmed mean – class of estimator of mean
 - A trimmed mean is computed by “trimming away” a certain percent of both the largest and the smallest set of values.
 - For ex: the 10% trimmed mean is found by eliminating the largest 10% and smallest 10% and computing the average of the remaining values.

Statistics and Data Analysis

Variability

- Sample Range and Sample standard deviation
 - Sample range: $X_{\max} - X_{\min}$
 - Sample standard deviation
 - Sample variance $s^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$, $n - 1$: degree freedom
 - Sample standard deviation: $\sqrt{s^2}$
 - The importance of measures of variability
 - Important population parameters:
 - population mean – related to sample standard deviation
 - population variance – sample variance .

Statistics and Data Analysis

Statistical modeling

- Statistical analysis
 - Estimation of parameters of a postulated model
 - Statistical model is random or stochastic rather than deterministic
 - Data Illustration - Graphical representation of a collection of data
 - Scatter plot
 - Histogram
 - Box plot

Statistics and Data Analysis

Statistical modeling

- Scatter plot
 - The sample means and variability are depicted nicely in the scatter plot.

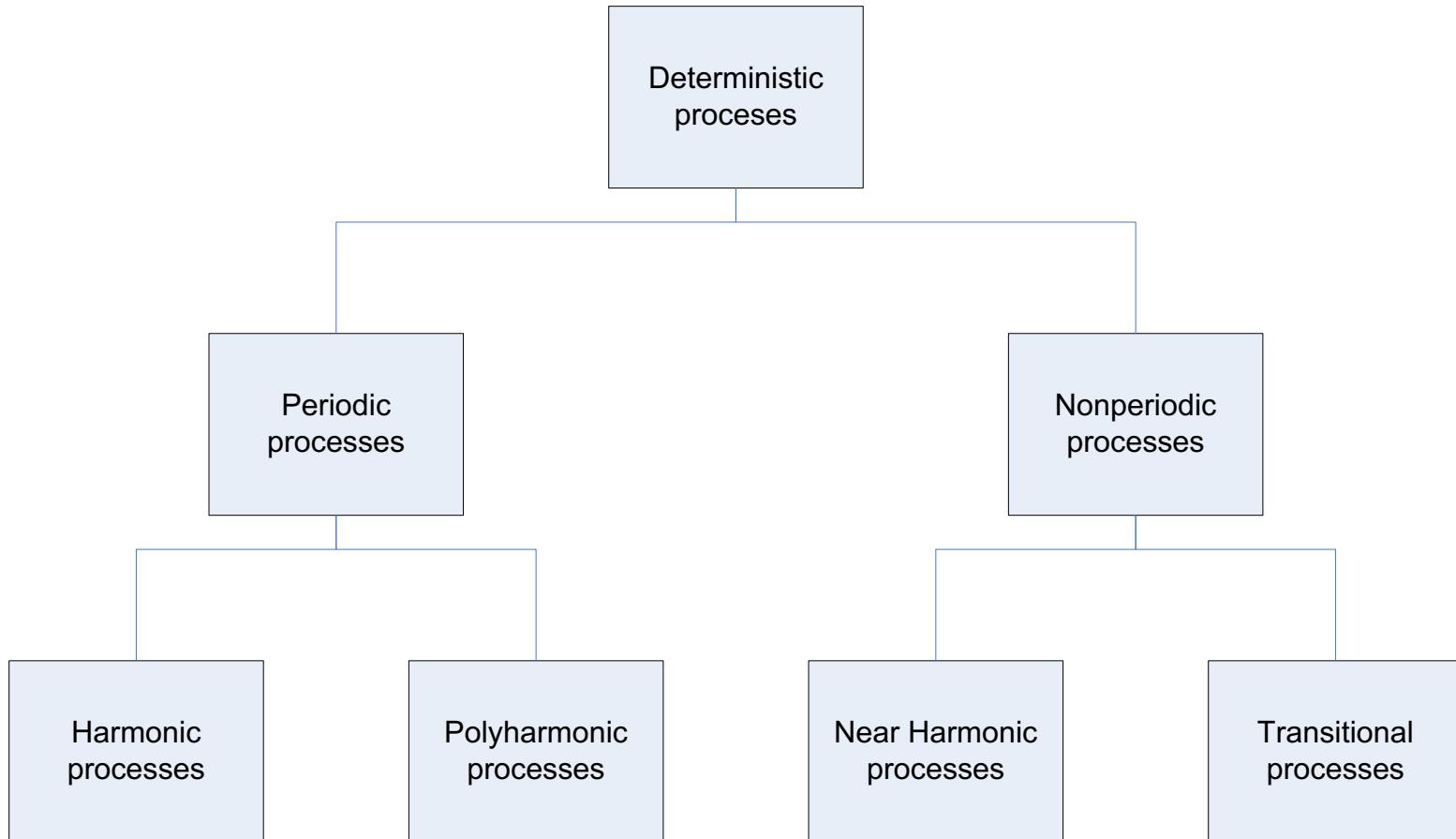
I. Models and Definitions

- Deterministic and stochastic processes
- Classification of deterministic processes
- Classification of stochastic processes
- Analysis of stochastic processes

1.1 Deterministic and stochastic processes

- Deterministic processes:
 - Definition: processes which can be described using explicit mathematic formulas
 - Example:
$$i(t) = A \cos(\omega_0 t + \varphi_0), \quad t \geq 0$$
 - Oscillation of current in linear RLC circuits
- Stochastic processes
 - Processes is random which is described using probability concepts and statistic characteristics.
 - Example:
 - Brownian motion.
 - Poisson process.

1.2. Classification of deterministic processes

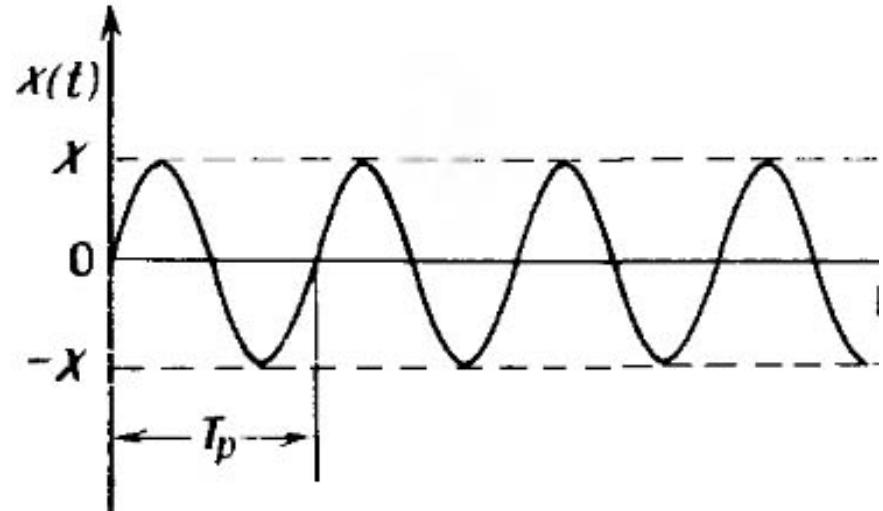
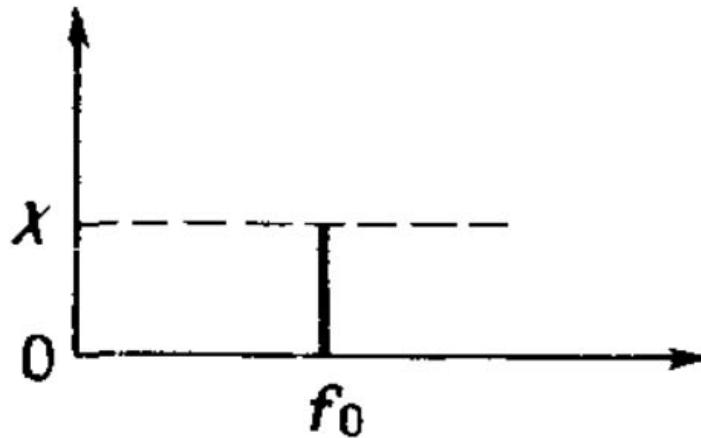


1.2. Classification of deterministic processes

- Periodic sinusoidal processes
 - Sinusoidal process – harmonic processes: the process is presented by:
$$x(t) = X \sin(2\pi f_0 t + \theta)$$
 - X: amplitude
 - f_0 : frequency, period $T_p = 1/f_0$
 - θ : phase shift
 - Example: sinusoidal current in the linear circuits.
 - In compact form: $x(t) = X \sin(2\pi f_0 t)$

1.2. Classification of deterministic processes

- Spectrum of sinusoidal processes
 - Line spectrum
- Harmonic processes are most simple deterministic processes



1.2. Classification of deterministic processes

□ Polyharmonic processes

- Process described by equation:

$$x(t) = x(t \pm nT_p)$$

- T_p : fundamental period
- $f = 1/T_p$: fundamental frequency
- Harmonic process: special case of polyharmonic processes.

1.2. Classification of deterministic processes

- Fourier series expansion of polyharmonic processes

$$x(t) = \frac{a_0}{2} + \sum_{n=1}^{\infty} (a_n \cos 2\pi n f_1 t + b_n \sin 2\pi n f_1 t)$$

$$f_1 = \frac{1}{T_p}$$

$$a_n = \frac{2}{T_p} \int_0^{T_p} x(t) \cos 2\pi n f_1 t dt, n = 0, 1, 2, \dots$$

$$b_n = \frac{2}{T_p} \int_0^{T_p} x(t) \sin 2\pi n f_1 t dt, n = 1, 2, 3, \dots$$

1.2. Classification of deterministic processes

- Other form of Fourier series expansion

$$x(t) = X_0 + \sum_{n=1}^{\infty} (X_n \cos 2\pi n f_1 t - \theta_n)$$

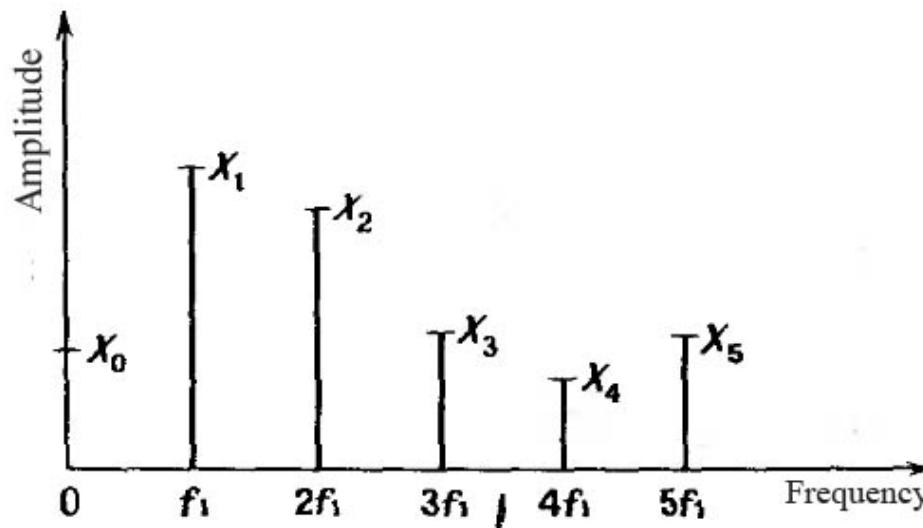
$$X_0 = \frac{a_0}{2}$$

$$X_n = \sqrt{a_n^2 + b_n^2}, n = 1, 2, 3, \dots$$

$$\theta_n = \operatorname{arctg} \left(\frac{a_n}{b_n} \right), n = 1, 2, 3, \dots$$

1.2. Classification of deterministic processes

- Spectrum of polyharmonic processes
 - Line spectrum



1.2. Classification of deterministic processes

- Near periodic processes
 - Periodic processes as usual are presented by series consisted of harmonic processes with commensurable frequencies
 - Conversely, any process, which can be presented as sum of two or more harmonic processes with commensurable frequency, is periodic.

1.2. Classification of deterministic processes

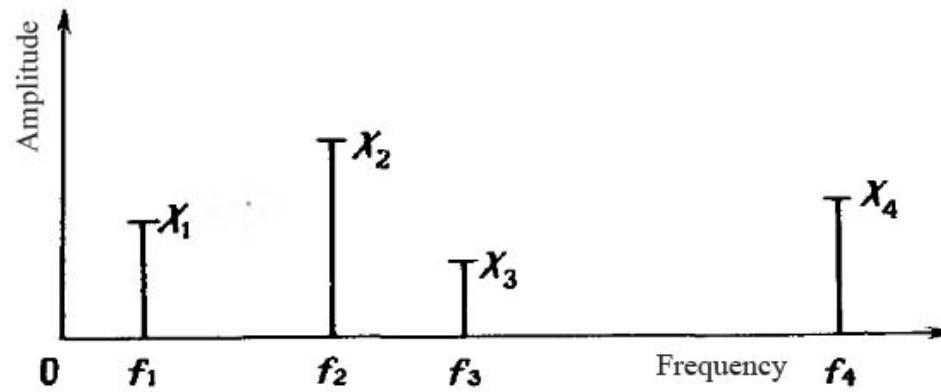
- Near periodic processes are presented by following equation:

$$x(t) = \sum_{n=1}^{\infty} (X_n \sin(2\pi n f_n t + \theta_n))$$

- Where f_n/f_m are rational in most of frequencies.

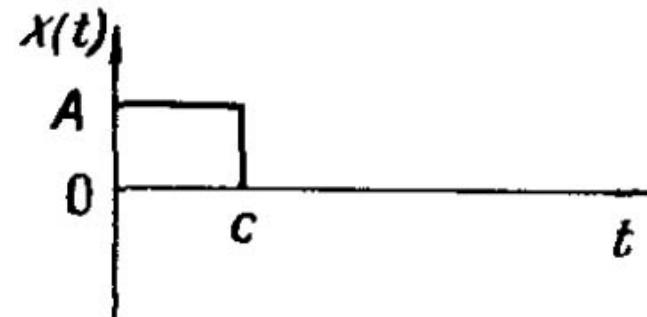
1.2. Classification of deterministic processes

- Near periodic processes has importal property:
 - If phase shift θ_n is omitted, near periodic process can be presented by line spectrum like polyharmonic process.
 - The distinction from periodic processes: frequency ratio of components is irrational



1.2. Classification of deterministic processes

- Transitional nonperiodic processes
 - In fact, almost all transitional processes are nonperiodic.
 - Example, oscillation of instant current in linear RLC circuit when current source is attached to the circuit



1.2. Classification of deterministic processes

- Distinction of transitional processes:
 - TPs can not be presented by line spectrum
 - Continuous spectrum:

$$X(f) = \int_{-\infty}^{\infty} x(t)e^{-j2\pi ft} dt$$

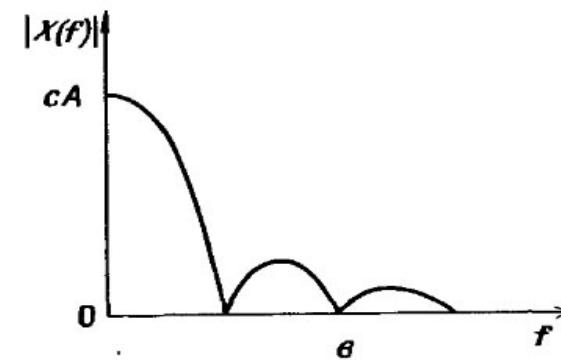
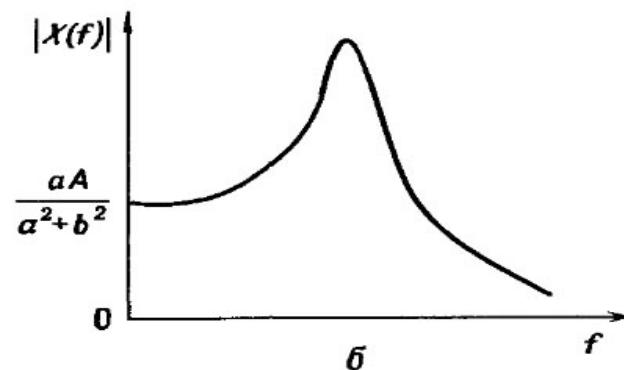
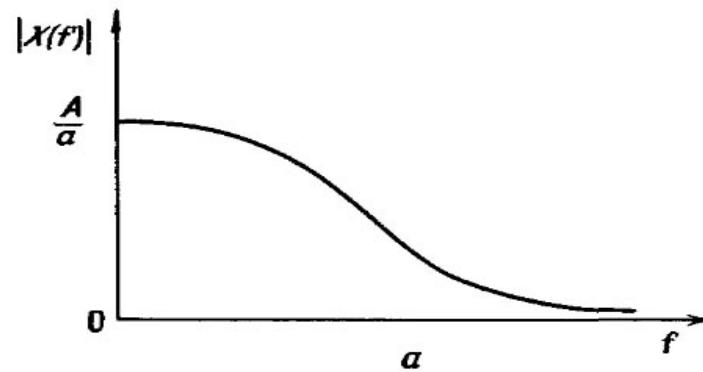
$$x(t) = 2\pi \int_{-\infty}^{\infty} X(f)e^{j2\pi ft} df$$

$$X(f) = |X(f)|e^{-j\theta(f)}$$

- $|X(f)|$ - module of $X(f)$
- $\theta(f)$ – argument of $X(f)$

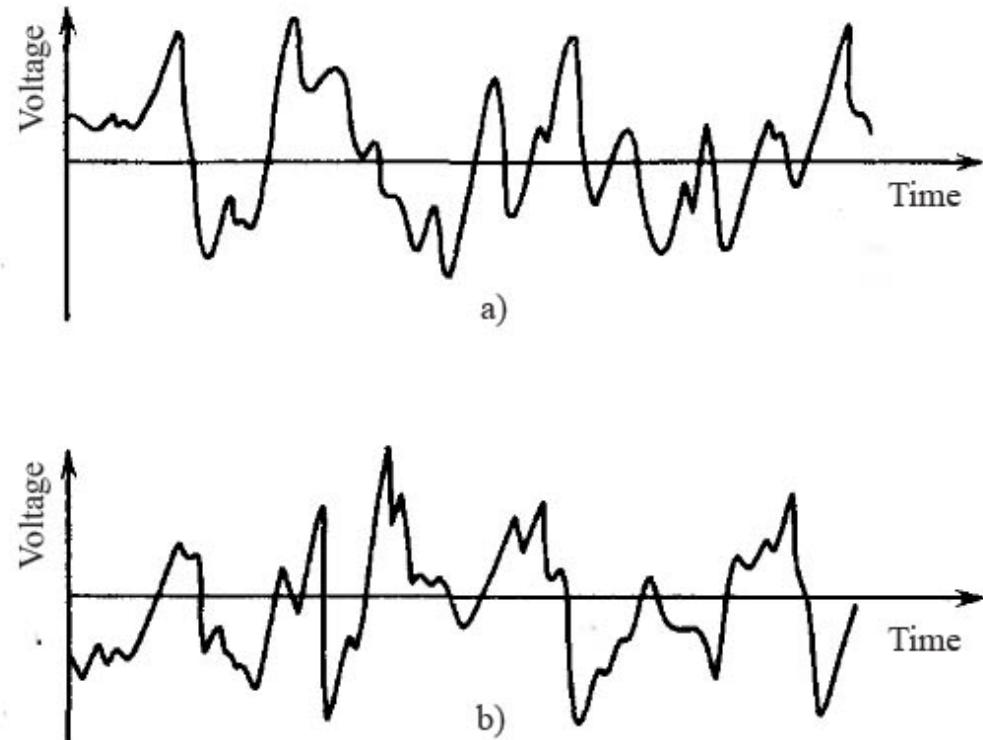
1.2. Classification of deterministic processes

- Spectrum of transitional processes



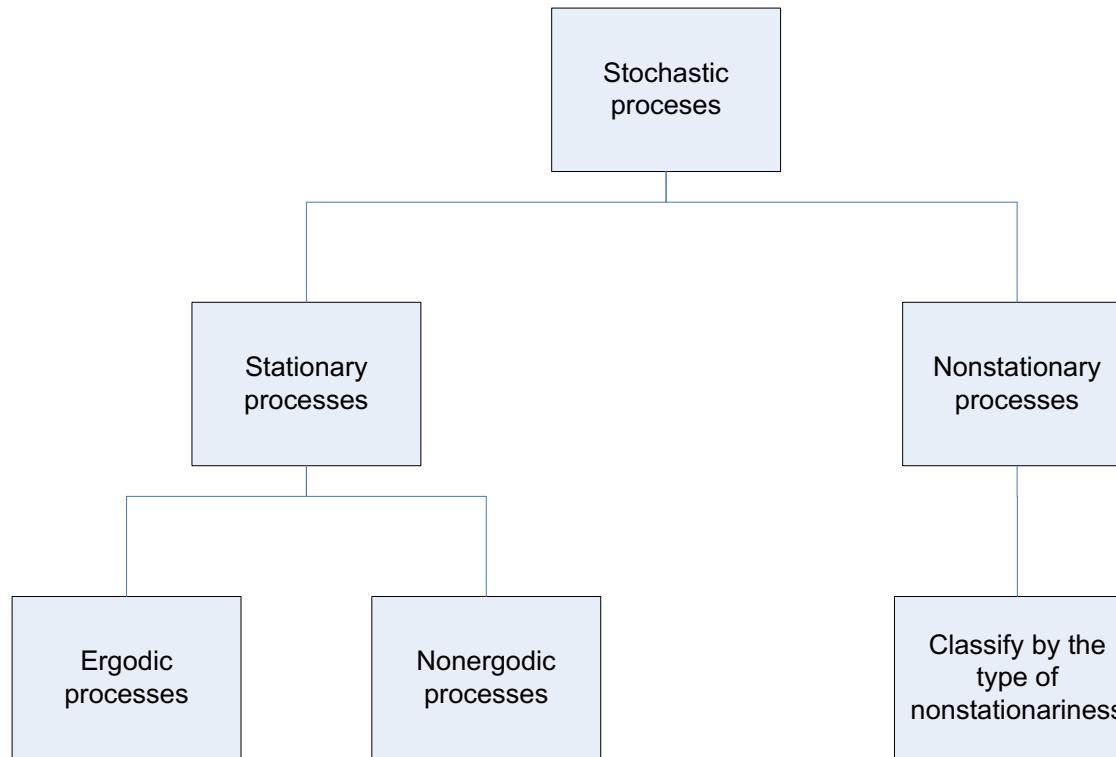
1.3. Classification of Stochastic processes

- Stochastic processes
 - Process, which describes random phenomenon, is called sample function
 - Ensemble of all possible sample function, which describes random phenomenon, is called stochastic process.



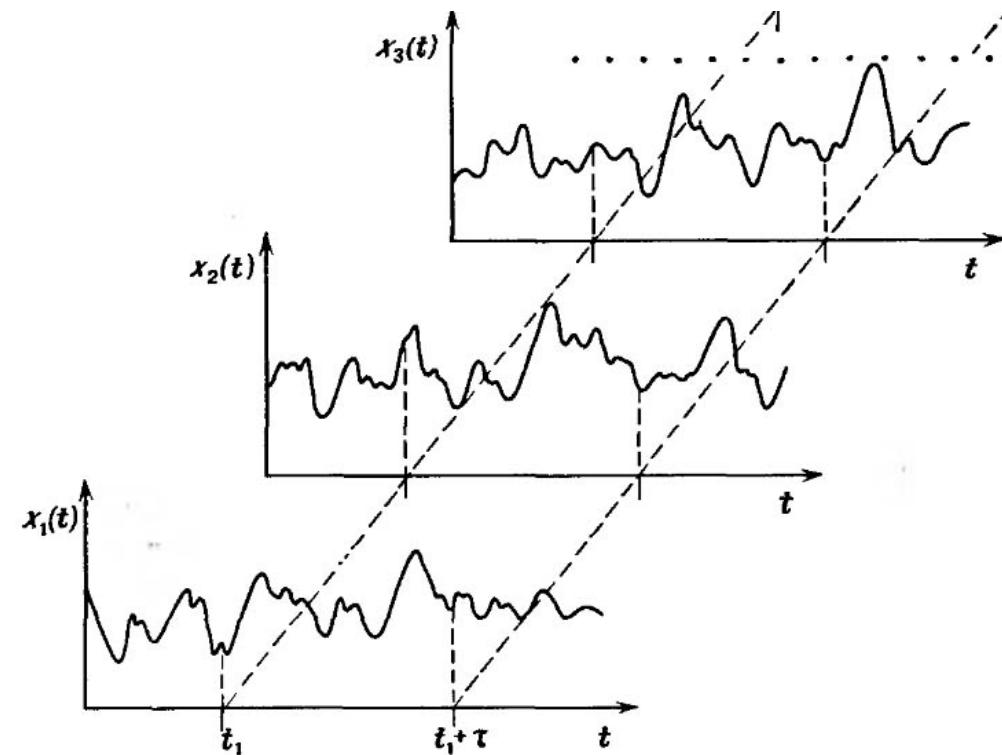
1.3. Classification of Stochastic processes

- Classification of stochastic processes



1.3. Classification of Stochastic processes

- Stationary stochastic process
 - Sample mean and covariance of random process



1.3. Classification of Stochastic processes

- Mean and covariance of random processes

$$\mu_x(t_1) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N x_k(t_1)$$

$$R_{xx}(t_1, t_1 + \tau) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N x_k(t_1)x_k(t_1 + \tau)$$

- If $\mu_x(t_1)$ and $R_{xx}(t_1, t_1 + \tau)$ does not depend time moment t_1 : process is weak stationary or stationary in wide sense.
- $\mu_x(t_1) = \mu_x$ and $R_{xx}(t_1, t_1 + \tau) = R_{xx}(\tau)$

1.3. Classification of Stochastic processes

- Strict sense stationary stochastic process
 - All of the moments and cross moments of stochastic process are time invariant.
 - In some cases, strict sense stationary is followed from wide sense stationary.
 - Wide sense and strict sense stationary.

1.3. Classification of Stochastic processes

- Ergodic stochastic processes
 - Time average and time covariance of kth sample function

$$\mu_x(k) = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T x_k(t) dt$$

$$R_{xx}(\tau, k) = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T x_k(t) x_k(t + \tau) dt$$

- If stochastic process {x(t)} stationary and characteristics $\mu_x(k)$ and $R_{xx}(\tau, k)$ of all its sample functions are equal, then process is called ergodic.

1.3. Classification of Stochastic processes

- Ergodic process has property: time average is equal ensemble average and the same is applicable to covariance.

$$\mu_x(k) = \mu_x$$

$$R_{xx}(\tau, k) = R_{xx}(\tau)$$

- Only stationary process can have ergodic property.

1.3. Classification of Stochastic processes

- Nonstationary stochastic processes
 - All stochastic processes, which is not meet fully stationary properties, are non-stationary.
 - Characteristics of nonstationary: depend on time, and can be calculated by averaging in separated time moment over ensemble of sample functions
 - Simplification of nonstationary processes.

1.3. Classification of Stochastic processes

- Stationarity of sample functions
 - Consider one sample function $x_k(t)$ of stochastic process $\{x(t)\}$
 - Average and covariance of $x_k(t)$ from time moment t_1 over time interval T are equal:

$$\mu_x(t_1, k) = \frac{1}{T} \int_{t_1}^{t_1+T} x_k(t) dt$$

$$R_{xx}(\tau, k) = \frac{1}{T} \int_{t_1}^{t_1+T} x_k(t) x_k(t + \tau) dt$$

1.3. Classification of Stochastic processes

- Nonstationary process: sample values are sounded with t_1 .
- Sample function of ergodic process is stationary
- Sample function of most of interesting nonstationary processes is nonstationary.

1.4. Analysis of stochastic processes

- Concepts
 - Realization of stochastic process can not be described by explicit math equations
 - Estimations of properties of these processes have to be done using methods of statistics
 - For specific applications, stochastic processes, which have properties satisfied determined relations, have important role in describing processes in reality.
 - It is important to take account in statistical errors related to parameter estimation and correlations between input and output processes of transformations.

1.4. Analysis of stochastic processes

- Basic characteristics of stochastic processes.
 - Important characteristics of stochastic processes
 - Mean and mean square
 - Probability density
 - Covariance function
 - Spectral density

1.4. Analysis of stochastic processes

- Mean μ_x and variance σ_x^2 of sample function of stochastic process describe center and square of value of data scatters
- Mean square ψ_x^2

$$\psi_x^2 = \sigma_x^2 + \mu_x^2$$

- Probability distribution density $p(x)$ of sample functions

$$p_x(x) = \frac{dF_x(x)}{dx}$$

1.4. Analysis of stochastic process

- $F_x(x) = P_x\{ \xi \leq x \}$: probability distribution of process $x(t)$.
- $F_x(x_2) - F_x(x_1) = P_x\{ x_2 \leq \xi \leq x_1 \}$
- Covariance function $R_{xx}(\tau)$ of stationary stochastic process gives measure of its value relation.

$$R_{xx}(\tau, t_1) = E\{x(t_1)x(t_1+\tau)\}$$

- Covariance function can be estimated by multiplying sample value in any time moment with value shifted to interval equal τ . Calculation is made for all required shifted interval.

1.4. Analysis of stochastic processes

- Spectral density $G_{xx}(f)$ of stationary stochastic processes gives speed of change of square average in dependency on frequency.
- For two sample functions of stochastic processes, joint statistic characteristics play important roles in its analysis.
 - Joint probability density function
 - Joint covariance function
 - Joint spectral density
 - Joint frequency characteristics
 - Coherent funtions

1.4. Analysis of stochastic processes

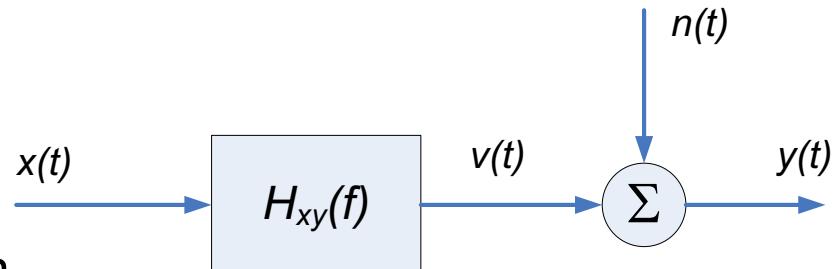
- Reasons of using Probability distribution function and probability density function:
 - Checking normality of the process
 - Identification of nonlinearity of the process
 - Analysis of extremal values of the process
- Applications of covariance functions
 - Identification of periodicity
 - Separate signal in noise
 - Measuring signal delay
 - Detection of noise signal source
 - Evaluation of signal transmission speed

1.4. Analysis of stochastic processes

- Typical applications of spectral density:
 - Definition of system properties by observations of input and output processes.
 - Prediction of output processes knowing input processes and system properties.
 - Identification of input processes by output processes and system properties.
 - Give dynamic data for testing.
 - Identification of signal and noise sources.
 - Optimal linear prediction and filtration.

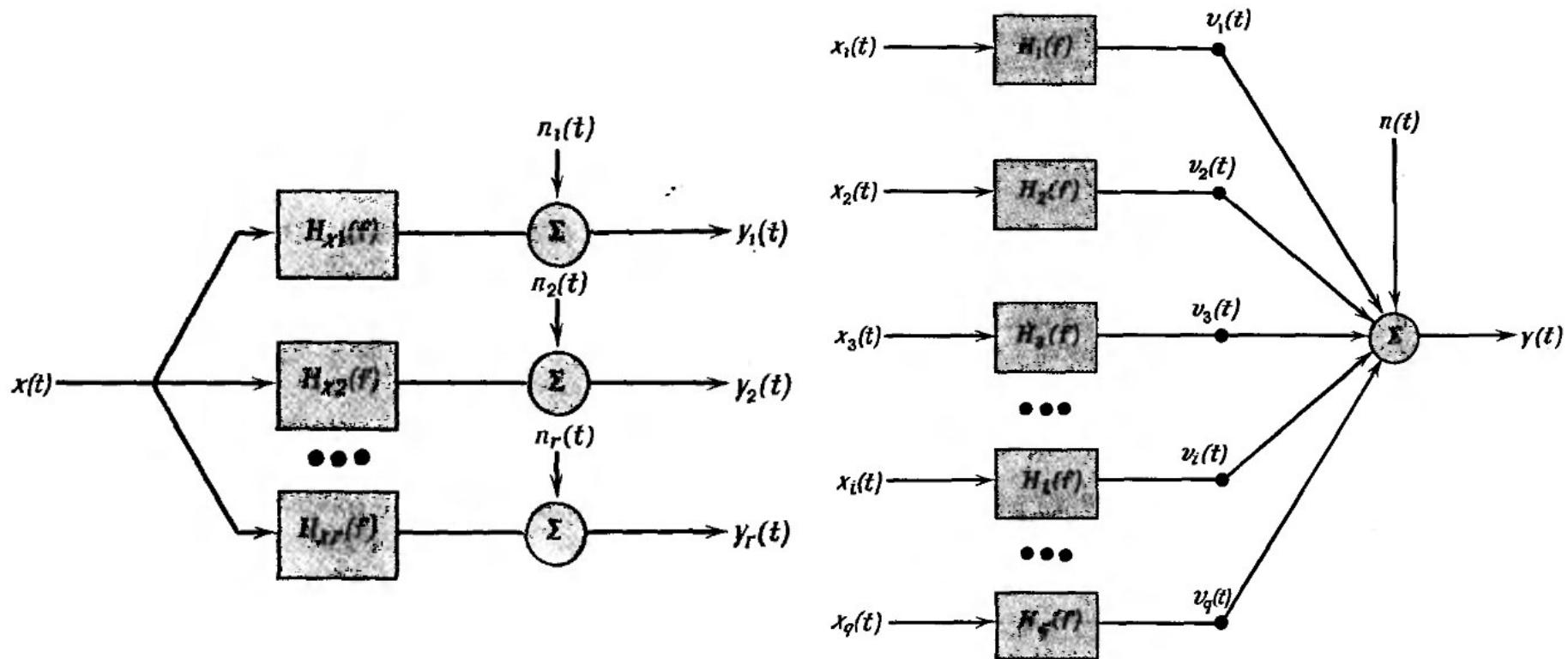
1.4. Analysis of stochastic processes

- Input and output relation
 - Models of input – output transformation:
 - Model with one input and one output.
 - Model with one input and many output.
 - Model with many input and one output.
 - Model with many input and many output.
 - Most simple model: model with one input and one output:
 - $n(t)$: noise signal
 - $y(t) = H_{xy}[x(t)] + n(t)$
 - If the system is linear
And time invariant, then
impact of the system to
input signal become convolution.



1.4. Analysis of stochastic processes

- Building complex system based on simple systems



1.4. Analysis of stochastic processes

- Characteristics of experimental errors
 - Estimation of characteristics of errors:
 - Estimation of parameter θ is denoted by $\hat{\theta}$
 - Value $\hat{\theta}$ is estimate of parameter θ in finite time interval or in finite number of sample values.
 - Suppose: repeated experiments for estimating $\hat{\theta}$
 - Expectation:

$$E\{\hat{\theta}\} = \frac{1}{N} \sum_{i=1}^N \hat{\theta}_i$$

- Unbiased estimation $b[\hat{\theta}] = E\{\hat{\theta}\} - \theta$
- Biased estimation
- Bias of estimation - systematic error of estimation.

1.4. Analysis of stochastic processes

- Variance of estimate
 - Definition:

$$Var[\hat{\theta}] = E\{(\hat{\theta} - E\{\hat{\theta}\})^2\}$$

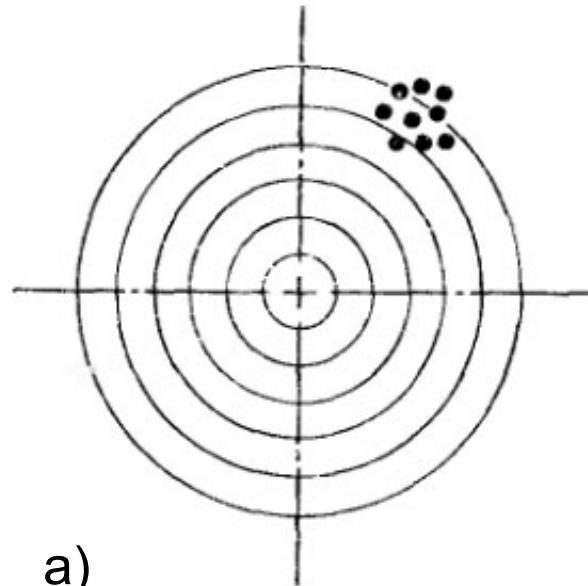
- Aggregate error of estimate:

$$E\{(\hat{\theta} - \theta)^2\} = Var[\hat{\theta}] + b[\hat{\theta}]^2$$

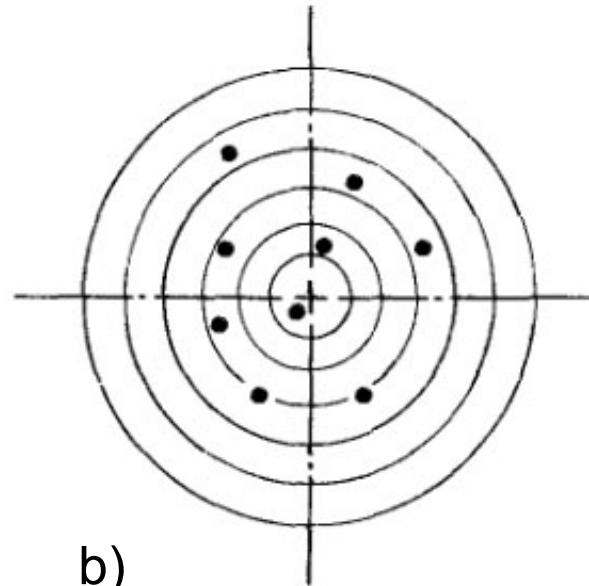
- If bias of estimation is equal 0 or approximately equal 0, then variance and aggregate error are concurred

1.4. Analysis of stochastic processes

- Example:



a)



b)

1.4. Analysis of stochastic processes

- Normalized average square root error

$$\varepsilon[\hat{\theta}] = \frac{\sqrt{E\{(\hat{\theta} - \theta)^2\}}}{\theta}$$

- If this error is as small as possible, then the estimate is close to true value.

1.4. Analysis of data

- Data analysis methods
 - Collecting data
 - Data gathering
 - Preparing data for preliminary rating
 - Analysis of stationary random processes
 - Filtering stochastic processes
 - Fourier analysis of random processes
 - Numerical methods for estimating probability density, covariance functions and other statistical characteristic.

Chapter resume

- Classes of deterministic processes
- Stochastic processes
 - Sample functions
 - Characteristics:
 - Probability distribution function.
 - Probability density function.
 - Mean and variance.
 - Covariance function.
- Data Analysis methods
 - Collecting data
 - Analysis of stationary stochastic processes
 - Filtering nonstationary processes
 - Fourier analysis
 - ...



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG





VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG





25
YEARS ANNIVERSARY
SOICT

VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

**Thank you for
your attentions!**

