# Probability And Statistics

**Prof. Dr. Serdar KORUKOĞLU**

## Ege University Department of Computer Engineering

# References

1. Jay L. Devore, **"Probability and Statistics for Engineering and Sciences",** Thomson, International Student Edition. Seventh Edition.

2. Ranold E. Walpole, Raymond H. Myers, Sharon L. Myers, Keying Ye. **"Probability and Statistics for Engineers and Scientists",** Pearson Educational International. $8^{th}$ Edition.

3. Douglas C. Montgomery, George C. Runger, **"Applied Statistics and Probability for Engineers",** John Wiley & Sons, Inc. $4^{th}$ Edition.

4. William Mendenhall, Robert J. Beaver, Barbara M. Beaver, **"Introduction to Probability and Statistics",** Brooks/Cole, $13^{th}$ Edition.

**The world is becoming more and more quantitative. <span style="color:red">Many professions depend on numerical measurements to make decisions in the face of uncertainty</span>. Statisticians use <span style="color:red">quantitative abilities</span>, <span style="color:red">statistical knowledge</span>, and communication skills to work on many challenging problems.**

**The use of Statistical methods in manufacturing development of computer software and other areas involves the gathering of information or <span style="color:red">scientific data.</span>**

*The word <span style="color:red">"statistics"</span> comes from the Latin word <span style="color:blue">"status"</span>, which means the state of phenomena; <span style="color:blue">stato – state, statista – statistic, state expert, statistika</span> – a certain sum of knowledge, information about the state.* <span style="color:red">From Wikipedia:</span>

*Statistics is the science of making effective use of numerical data relating to groups of individuals or experiments. <u>It deals with all aspects of this, including not only the collection, analysis and interpretation of such data, but also the planning of the collection of data, in terms of the design of surveys and experiments</u>.*

*A <span style="color:red">statistician</span> is someone who is particularly versed in the ways of thinking necessary for the successful application of statistical analysis. Often such people have gained this experience after starting work in any of a list of fields of application of statistics. There is also a discipline called mathematical statistics, which is concerned with the theoretical basis of the subject.*

*The word <span style="color:red">statistics can either be singular or plural</span>. In its <span style="color:blue">plural form</span>, statistics refers to the mathematical science discussed in this article. In its <span style="color:blue">singular form</span>, statistics is the plural of the word statistic, which <span style="color:red">refers to a quantity</span> (such as a mean) calculated from a set of data.*

**Statistics teaches us how to make intelligent judgments and informed decisions in the presence of uncertainity or variation.**

## Computer Software

- **SAS**

- **MINITAB**

- **SPSS**

- **IMSL Library**

- **NAG Library**

- **R, Open Source Software** http://www.r-project.org/

**More than 4,000 universities** worldwide rely on Minitab to teach statistical concepts and simplify learning. We help teachers teach and students understand statistics.

And because thousands of companies rely on our software for data analysis and quality improvement, students who learn with Minitab also get experience with real-world business tools.



- **Beginning in the 1980s and continuing into the twenty-first century; an inordinate amount of attention has been focused on improvement of quality.**

- **Much has been said and written about the Japanese "Industrial Miracle" which began in the middle of the twentieth century. The Japanese were able to succeed where we and other countries had failed namely, to create an atmosphere that allows the production of high-quality products. Much of the success of the Japanese has been attributed to the use of statistical methods and statistical thinking among management personnel.**

# The R Project for Statistical Computing



R is a language and environment for statistical computing and graphics. It is a GNU project which is similar to the S language and environment which was developed at Bell Laboratories (formerly AT&T, now Lucent Technologies) by John Chambers and colleagues. R can be considered as a different implementation of S. There are some important differences, but much code written for S runs unaltered under R.

R provides a wide variety of statistical (linear and nonlinear modeling, classical statistical tests, time-series analysis, classification, clustering, ...) and graphical techniques, and is highly extensible. The S language is often the vehicle of choice for research in statistical methodology, and R provides an Open Source route to participation in that activity.

One of R's strengths is the ease with which well-designed publication-quality plots can be produced, including mathematical symbols and formulae where needed. Great care has been taken over the defaults for the minor design choices in graphics, but the user retains full control.

R is available as Free Software under the terms of the Free Software Foundation's GNU General Public License in source code form. It compiles and

runs on a wide variety of UNIX platforms and similar systems (including FreeBSD and Linux), Windows and MacOS.

**The R environment**

R is an integrated suite of software facilities for data manipulation, calculation and graphical display. It includes

- an effective data handling and storage facility,
- a suite of operators for calculations on arrays, in particular matrices,
- a large, coherent, integrated collection of intermediate tools for data analysis,
- graphical facilities for data analysis and display either on-screen or on hardcopy, and
- a well-developed, simple and effective programming language which includes conditionals, loops, user-defined recursive functions and input and output facilities.

The term "environment" is intended to characterize it as a fully planned and coherent system, rather than an incremental accretion of very specific and inflexible tools, as is frequently the case with other data analysis software.

R, like S, is designed around a true computer language, and it allows users to add additional functionality by defining new functions. Much of the system is itself written in the R dialect of S, which makes it easy for users to follow the algorithmic choices made. For computationally-intensive tasks, C, C++ and Fortran code can be linked and called at run time. Advanced users can write C code to manipulate R objects directly.

Many users think of R as a statistics system. We prefer to think of it of an environment within which statistical techniques are implemented. R can be extended (easily) via *packages*. There are about eight packages supplied with the R distribution and many more are available through the CRAN family of Internet sites covering a very wide range of modern statistics.

R has its own LaTeX-like documentation format, which is used to supply comprehensive documentation, both on-line in a number of formats and in hardcopy.

- **ELEMENTARY STATISTICS**
- **INTERMEDIATE STATISTICAL CONCEPTS**
- **HONORS ELEMENTARY STATISTICS**
- **STATISTICS FOR SCIENTISTS & ENGINEERS**
- **STATISTICAL COMPUTING AND GRAPHICAL ANALYSIS**
- **STATISTICS FOR RESEARCHERS**
- **METHODS FOR DATA ANALYSIS**
- **PROBABILITY**
- **MATHEMATICAL STATISTICS**
- **NONPARAMETRIC STATISTICS**
- **INTRODUCTION TO TIME SERIES ANALYSIS**
- **INTRODUCTION TO APPLIED MULTIVARIATE ANALYSIS**
- **INTRODUCTION TO CATEGORICAL DATA ANALYSIS**
- **SAMPLING**
- **MIXED EFFECTS MODELS**
- **INTERMEDIATE PROBABILITY & STATISTICS**
- **INTERMEDIATE MATHEMATICAL STATISTICS**
- **LINEAR MODELS**
- **ADVANCED REGRESSION ANALYSIS**
- **TOPICS IN APPLIED STATISTICS**
- **STOCHASTIC PROCESSES**
- **BIOSTATISTICS**
- **EXPERIMENTAL DESIGN**
- **STATISTICAL QUALITY CONTROL**
- **BAYESIAN DATA ANALYSIS**
- **SPATIAL DATA ANALYSIS**
- **INTRODUCTION TO TIME SERIES ANALYSIS**
- **MULTIVARIATE ANALYSIS**
- **GENERALIZED LINEAR MODELS**
- **ADVANCED MATHEMATICAL STATISTICS**
- **RESPONSE SURFACE METHODOLOGY**
- **DATA MINING**

# Engineers must know how to

- **Efficiently plan experiments,**

- **Collect data ,**

- **Analyze and the interpret the data, and**

- **Understand how the observed data are related to the model they have proposed for the problem under study.**

# Engineering Method

The steps is the engineering method are as follows:

1. Develop a clear and concise description of the problem.
2. Identify, at least tentatively, the important factors that affect this problem or that may play a role in its solution.
3. Propose a model for the problem, using scientific or engineering knowledge of the phenomenon being studied. State any limitations or assumptions of the model.
4. Conduct appropriate experiments and collect data to test or validate the model or conclusions made in steps 2 and 3.
5. Refine the model to assist in developing a solution to the problem.
6. Manipulate the model to assist in developing a solution to the problem.
7. Conduct an appropriate experiment to confirm that the proposed solution to the problem is both affective and efficient.
8. Draw conclusions or make recommendations based on the problem solution.

```
┌─────────────────────┐
│   Develop a clear   │
│     description     │
└─────────────────────┘
          │
          ▼
┌─────────────────────┐      ┌─────────────────────┐
│ Identify the important │───▶│ Conduct experiments │
│       factors       │      └─────────────────────┘
└─────────────────────┘                 │
          │                             │
          ▼                             │
┌─────────────────────┐◀───────────────┘
│  Propose or refine a │
│        model        │
└─────────────────────┘
          │
          ▼
┌─────────────────────┐
│ Manipulate the model │
└─────────────────────┘
          │
          ▼
┌─────────────────────┐
│ Confirm the solution │
└─────────────────────┘
          │
          ▼
┌─────────────────────┐
│   Conclusions and   │
│   recommendations   │
└─────────────────────┘
```

# The Engineering Method

# The field of statistics deals with the

- **Collection,**

- **Presentation,**

- **Analysis, and**

- **Use of data to make decisions,**

- **Solve problems, and**

- **Design products and process.**

Because many aspects of engineering practice involve working with data, obviously some knowledge of statistics is important to an engineer.

Specifically, statistical techniques can be powerful aid in
- Designing new products, and systems,
- Improving existing designs,
- Designing, developing, and improving production processes.

# Variability

Statistical methods are used to help us describe and understand **variability**. By **variability** we mean that successive observations of a system or phenomenon do not product exactly the same result. Different factors represent potential **sources of variability** in the system.

Statistics gives us a framework for describing the **variability** and for learning about which <u>potential sources of</u> **variability** are the <u>most important.</u>

<u>**Statistics teaches us** how to make intelligent judgments **and informed decisions in the presence of** **uncertainity or variation**</u>.

# Statistical Inference



Often, physical laws (such as Ohm's Law) are applied to design products and processes.

But it is also important to reason from a specific set of measurement to more general cases. This reasoning is from a **sample** to a **population**. The reasoning is referred to as **statistical inference**.

# The Population and the sample

**Definition** A **population** is the set of all measurements of interest to the investigator.

**Definition** A **sample** is a subset of measurements selected from the population of interest.

In the language of statistics, one of the **most basic concepts i**s **sampling**. In most statistical problems, a specified number of measurements or data – a **sample** – is drawn from a much larger body of measurements, called the **population**.



**We try to describe or predict the behavior of the population on the basis of information obtained from a representative sample from that population.**

# Simple Random Sampling

**Simple Random Sampling** implies that any particular sample of a specified sample size has the same **chance** of being selected as any other sample of the same size.

Simple random sampling is not always appropriate.

The sampling units are not homogeneous and naturally divide themselves into nonoverlapping groups that are homogeneous. These groups are called **strata** and a procedure *called stratified random sampling* involves random selection of a sample within each **stratum**.

# Descriptive and Inferential Statistics

- **Descriptive Statistics** consist of procedures used to **summarize** and **describe** the important **characteristic** of a set of measurement. Computer-generated graphics and numerical summaries are commonplace in our everyday communications.

- **Inferential Statistics** consist of procedures used to make inferences about population characteristics from information contained in a sample drawn from this population.

**The objective of inferential statistics is to make inferences (that is draw conclusions, make predictions, make decisions) about the characteristics of a population from information contained in a sample.**

Probability

Population                    Sample

Statistical Inference

# Collecting Engineering Data

- **Retrospective Study**

A retrospective study would use either all or a sample of the <u>historical process</u> data archived over some period of time.

- **Observational  Study**

 In an observational study, the engineer <u>observes the process or population</u>, disturbing it as little as possible, and records the quantities of interest.

- **Designed Experiments**

In a designed experiment the <u>engineer makes deliberate or purposeful changes in the controllable variables of the system or process</u>, observes the resulting system output data, and then makes an inference or decision about which variables are responsible for the observed changes in output performance.

# Markov Chains

## Repair facility problem

A machine is either "working" or in the "repair center." If it is working today, then there is a 95% chance that it will be working tomorrow. It if is in the repair center today, then there is a 40% chance that it will be working tomorrow. We will be interested in questions like "what fraction of time does my machine spend in the repair shop?"

**Question:** Describe the DTMC for the above problem.

**Answer:**

- There are 2 states: "Working" and "Broken," where "Broken" denotes that the machine is in repair.

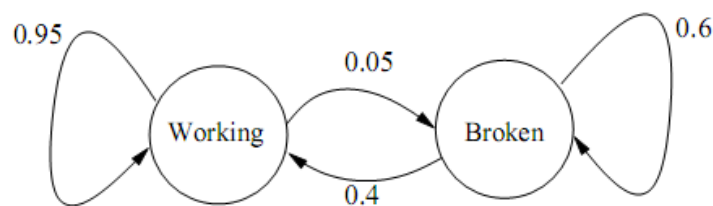- The Markov chain diagram is shown in Figure 1.



Figure 1: *Markov chain for simplest repair facility problem.*

- The Transition Probability Matrix is $\mathbf{P} = \begin{bmatrix} .95 & .05 \\ .4 & .6 \end{bmatrix}$
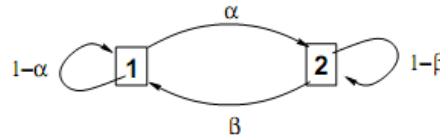
# A mouse in a cage

A mouse is in a cage with two cells, 1 and 2, containing fresh and stinky cheese, respectively. A mouse lives in the cage. A scientist's job is to record the position of the mouse every minute. When the mouse is in cell 1 at time $n$ (minutes) then, at time $n+1$ it is either still in 1 or has moved to 2.



Statistical observations lead the scientist to believe that the mouse moves from cell 1 to cell 2 with probability $\alpha = 0.05$; it does so, regardless of where it was at earlier times. Similarly, it moves from 2 to 1 with probability $\beta = 0.99$.

We can summarise this information by the TRANSITION DIAGRAM:



Another way to summarise the information is by the $2 \times 2$ TRANSITION PROBABILITY MATRIX

$$P = \begin{pmatrix} 1-\alpha & \alpha \\ \beta & 1-\beta \end{pmatrix} = \begin{pmatrix} 0.95 & 0.05 \\ 0.99 & 0.01 \end{pmatrix}$$

Questions of interest:
1. How long does it take for the mouse, on the average, to move from cell 1 to cell 2?
2. How often is the mouse in room 1 ?

Question 1 has an easy, intuitive, answer: Since the mouse really tosses a coin to decide whether to move or stay in a cell, the first time that the mouse will move from 1 to 2 will have mean $1/\alpha = 1/0.05 \approx 20$ minutes. (This is the mean of the binomial distribution with parameter $\alpha$.)

# Birth-Death Process

- General case:



$1 - p_0$   $1 - p_1 - q_1$   $1 - q_N$

$p_0$   $p_1$

$q_1$   $q_2$   $q_N$

$p_i$

$q_{i+1}$

- Locally, we have:

- Balance equations:   $\pi_i p_i = \pi_{i+1} q_{i+1}$

- Why? (More powerful, e.g. queues, etc.)

# M/M/1 Queue (1)

- Poisson **arrivals** with rate $\lambda$
- Exponential **service time** with rate $\mu$
- $\boxed{m = 1 \ \textbf{server}}$
- Maximum **capacity** of the system $= N$

- Discrete time intervals of (small) length $\delta$ :



- Balance equations: $\qquad \lambda \pi_{i-1} = \mu \pi_i \quad i \le N$
- Identical solution to the random walk problem.


# The Phone Company Problem (1)

- Poisson arrivals (**calls**) with rate $\lambda$
- Exponential service time (**call duration**), rate $\mu$
- $\boxed{m = N \ \text{servers (\textbf{number of lines})}}$
- Maximum capacity of the system $= N$

- Discrete time intervals of (small) length $\delta$ :



- Balance equations: $\qquad \lambda \pi_{i-1} = i\mu \pi_i$

- Solve to get: $\qquad \pi_i = \pi_0 \dfrac{\lambda^i}{\mu^i i!} \qquad \pi_0 = 1 / \displaystyle\sum_{i=0}^{N} \dfrac{\lambda^i}{\mu^i i!}$

# Designed experiment

A series of runs, or tests, in which you purposefully make changes to input variables simultaneously and observe the responses. A designed experiment is an efficient approach for improving a process because you can change more than one factor at a time to quickly obtain meaningful results and draw conclusions about how factors interact to affect the response.

**Minitab offers the following experimental designs:**

· **Factorial designs**

· **Response surface designs**

· **Mixture designs**

· **Taguchi designs**

# Variables and Data

- A **variable** is a characteristic that changes or varies over time and/or for different individuals or objects under consideration.

- An **experimental unit** is the individual or object on which a variable is measured.

- **Univariate data** result when a single variable is measured on a single experimental unit.

- **Bivariate data** result when two variables are measured on a single experimental unit.

- **Multivariate data** result when more than two variables are measured.

# Types of Variables

**Variables can be classified into one of two categories:**

## Qualitative or quantitative

**Definition Qualitative variables** measure a quality or characteristic on each experimental unit. **Quantitative variables** measure a <u>numerical quantity or amount</u> on each experimental unit.

## Examples (Qualitative)

- **Political affiliation: Republican, Democrat, Independent.**
- **Taste Ranking: excellent, good, fair, poor.**
- **Color: Brown, yellow, red, green, blue.**

## Examples (Quantitative)

- **x=Prime interest rate.**
- **x= Number of passengers on a flight from Izmir to Ankara.**
- **x= Weight of a package ready to shipped.**

```
                          ┌──────────────┐
                          │     Data     │
                          └──────────────┘
                          │
              ┌───────────┴────────────┐
              ▼                        ▼
      ┌──────────────┐         ┌──────────────┐
      │  Qualitative │         │ Quantitative │
      └──────────────┘         └──────────────┘
                                │
                        ┌───────┴────────┐
                        ▼                ▼
                ┌──────────────┐  ┌──────────────┐
                │   Discrete   │  │  Continuous  │
                └──────────────┘  └──────────────┘
```

A **discrete variable** can assume only a <u>finite or countable number</u> of values. A **continuous variable** can assume the <u>infinitely many values</u> corresponding to the points on a line interval.

**Random variable** when the values obtained arise as a result of <u>chance factors</u>, the variable is called a random variable.

**Discrete Random Variable: A Discrete Random Variable** is characterized by gaps or interruptions in the values that it can assume.

**Continuous Random Variable: A continuous random variable** does not possess the gaps or interruptions characteristic of a discrete random variable. Continuous random variables include the various measurements that can be made on individuals such as height and weight. No matter how close together the observed heights of two people, for example, we can, theoretically, find another person whose height falls somewhere in between.

# Data Type (Minitab)

**Refers to the different kinds of data recognized by Minitab: numeric, date/time, and text. Most Minitab analyses require data of specific types.**

**Data type is denoted in the Data window. As you can see in the illustration below, columns containing:**

- **Date/time data are marked with a D**
- **Text are marked with a T**
- **Numeric data are unmarked.**

| C1 | C2-D | C3-T | C4 |
|---|---|---|---|
| Numeric | Date/Time | Text | Currency |
| 4 | 2004/07/26 | Red | 21.75 |
| 7 | 2004/08/11 | Blue | 18.25 |
| 12 | 2004/12/02 | Green | 5.85 |
| 10 | 2004/12/26 | Yellow | 3.25 |
| 8 | 2005/01/01 | Orange | 24.50 |
| 5 | 2005/01/15 | Brown | 32.16 |

**You can <u>change the data type of a column</u>. If, for example, you enter or import data incorrectly and format a numeric column as text, the column can't be used in any analysis requiring numeric data until its data type is changed.**

# Measurement Scales

**Measurement:** This may be defined as the assignment of numbers to objects or events according to a set of rules. The various measurement scales result from the fact that measurement may be carried out under different sets of rules.

**The Nominal Scale:** The lowest measurement scale is the nominal scale. As the name implies it consists of "naming" observations or classifying them into **various mutually exclusive categories**.

- **Male-female**
- **Child-adult**
- **Under 20 years of age-20 and over**

**The Ordinal Scale:** Whenever observations are not only different from category to category, but can be ranked according to some criterion, they are said to be measured on an ordinal scale. Individuals may be classified according to socioeconomic status as low, medium or high. The intelligence of children may be above average, average or below average. The members of any one category are all considered equal, but the members of one category are all considered lower, worse or smaller than those in another category, which in turn bears a similar relationship to another category.

**The Interval Scale:** The interval scale is a more sophisticated scale than the nominal or ordinal in that with this scale it is not only possible to order measurements, but also the distance between any two measurements is known. We know, say that the difference between a measurement of 20 and a measurement of 30 is equal to the difference between measurements of 30 and 40. The ability to do this implies the use of a unit distance and a zero point, both of which arbitrary. The selected zero is not a true zero in that it does not indicate a total absence of the quantity being measured. (Temperature,Dates)

(100° is 10° warmer than 90°), we cannot multiply values or create ratios (100° is not twice as warm as 50°).

**The Ratio Scale:** The highest level of measurement is the ratio scale. This scale is characterized by the fact that equality of ratios as well as equality of intervals may be determined. Fundamental to the ratio scale is a true zero point. The measurement of such familiar traits as height, weight, and length makes use of the ratio scale.