# *Statistics and Probability*

*Hafara Firdausi, M.Kom.*

Department of Information Technology
Faculty of Electrical and Intelligent Informatics
Institut Teknologi Sepuluh Nopember

## 02. Concepts in Statistics
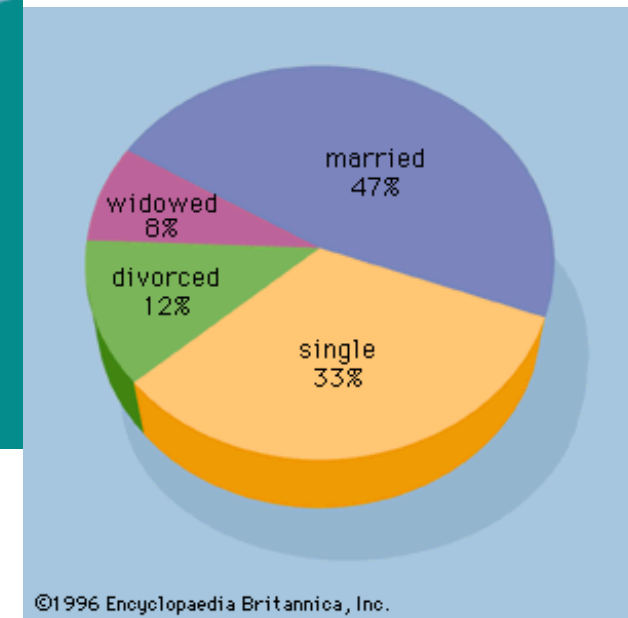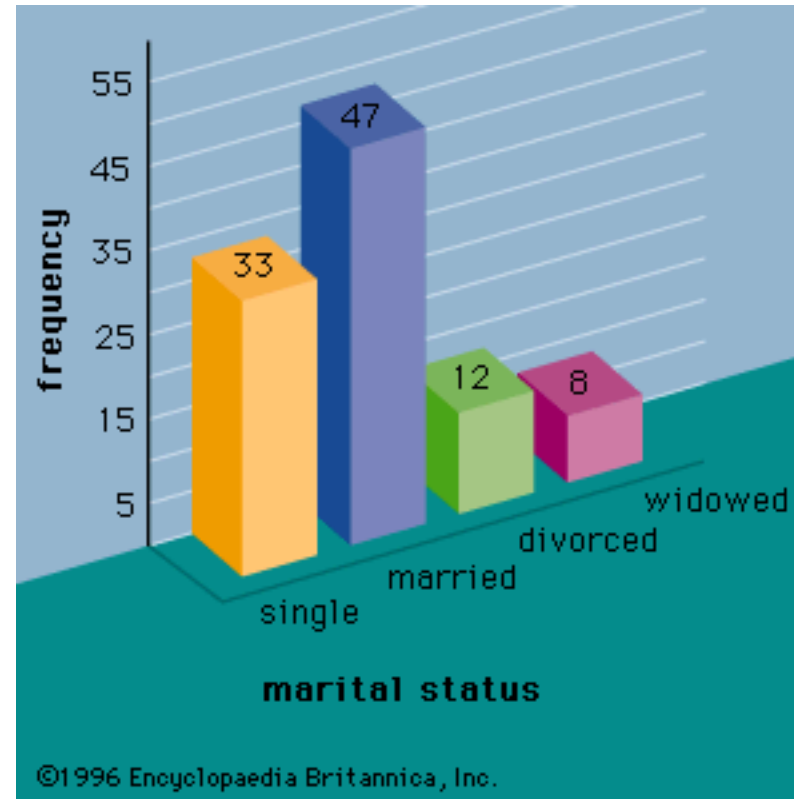
*Konsep Statistika*

# Outline

1. What is Statistics?

2. Sample and Population

3. Summary Measures
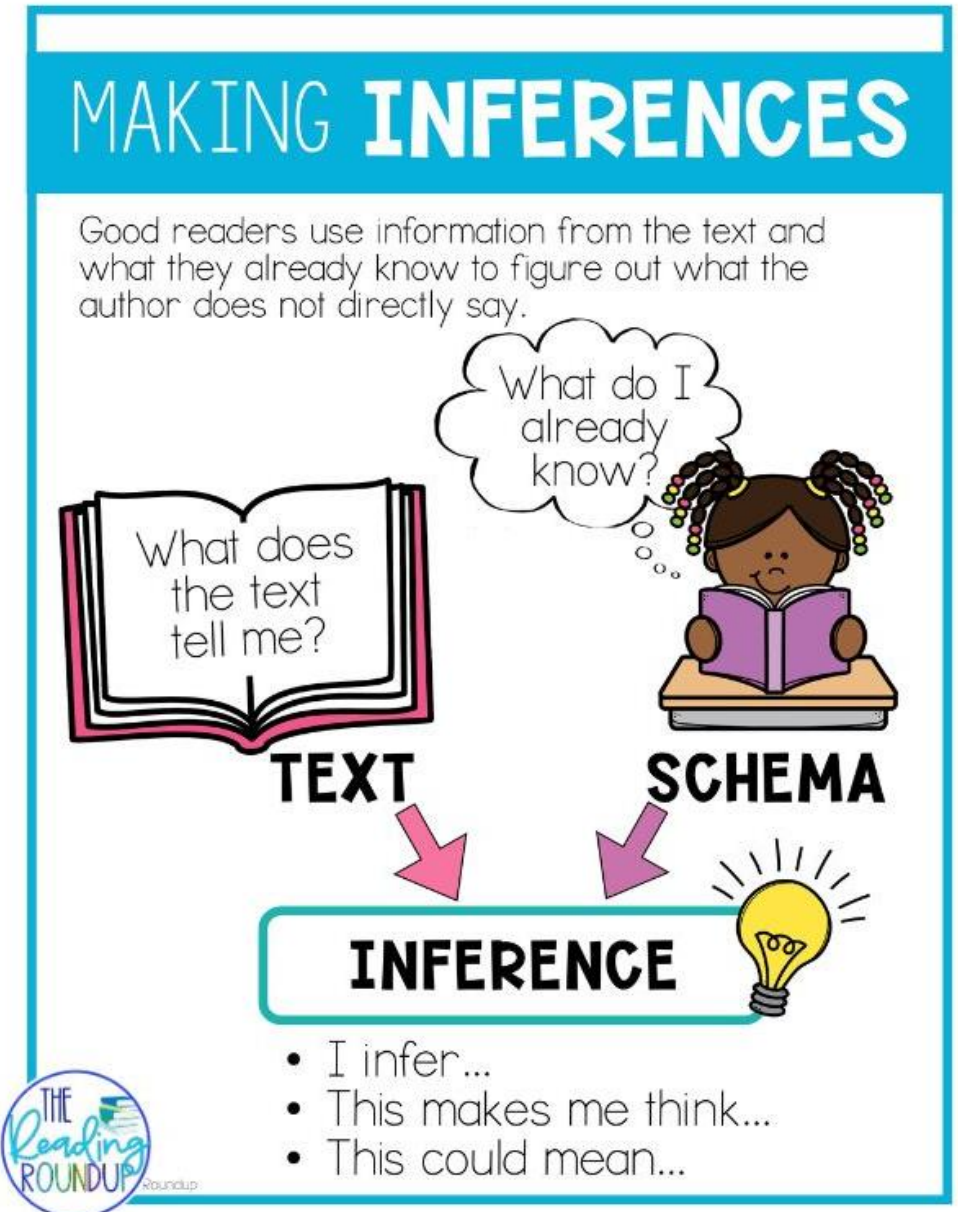
# What is Statistics?

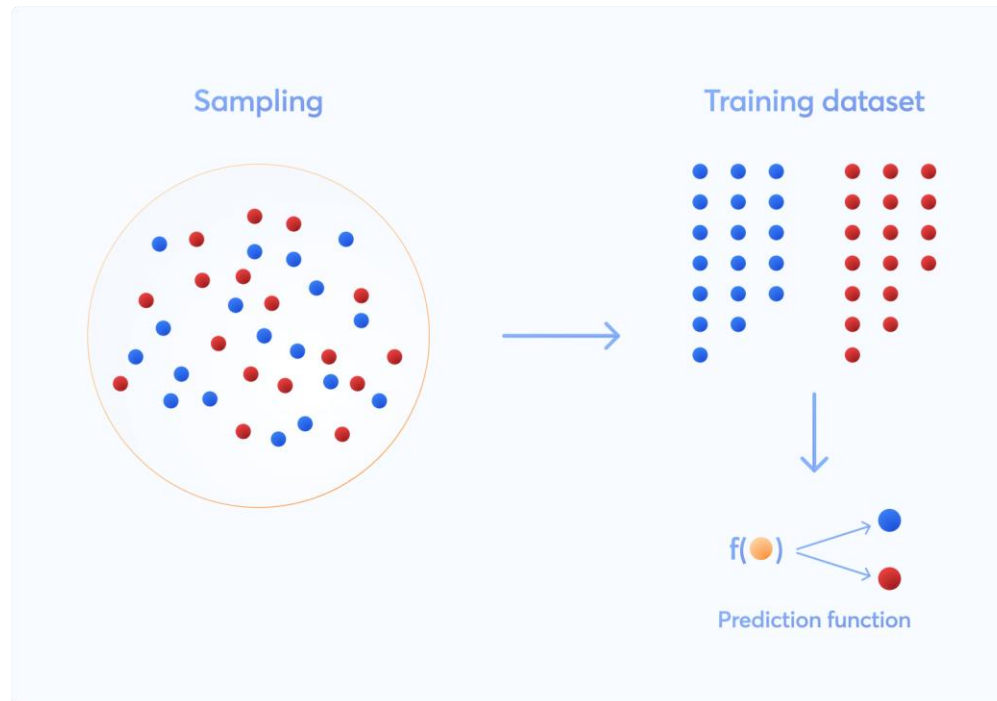# *What is Statistics?*

o Statistics is a **branch of mathematics** that consists of a set of **analytical techniques** that can be **applied to data** to help us make judgments and decisions in problems involving **uncertainty**.

o Statistics is a scientific discipline consisting of procedures for
   o **collecting**,
   o **describing**,
   o **analyzing**, and
   o **interpreting** numerical data

# *What is Statistics?*

**Main Objectives:**

To provide a set of procedures that enables us to **make inferences, predictions, and decisions** about **characteristics of a population of data** based on the information obtained from only a part of the population (**sample**).

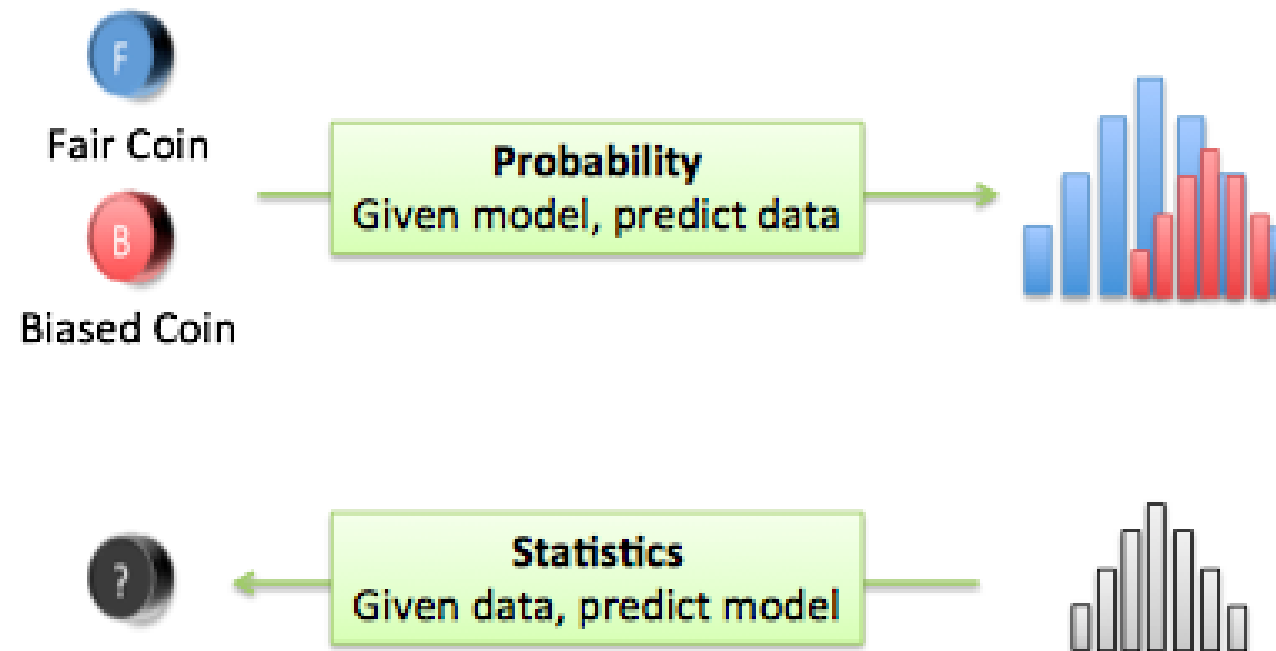# What is the Relation between Statistics and Probability?

"**Probability** provides a mathematical framework for **measuring uncertainty**"

"**Statistics** uses data to draw **conclusions** that are supported by that probability"



Fair Coin
Biased Coin

**Probability**
Given model, predict data

**Statistics**
Given data, predict model

*INSTITUT TEKNOLOGI SEPULUH NOPEMBER, Surabaya - Indonesia*

# *Case Study*

**Analisis faktor paling berpengaruh dalam penjualan Smartphone**

**D** (database)

→ **Statistika**
- Analisis regresi linear berganda untuk menentukan faktor yang berpengaruh dalam penjualan smartphone

→ **Probabilitas**
- Menghitung koefisien setiap variabel untuk melihat seberapa signifikan pengaruh variable tersebut terhadap penjualan

→ **Kesimpulan**

**Data penjualan smartphone**

Variabel:
- Resolusi Kamera
- RAM
- Harga
- Rating Pelanggan
- Jumlah Promosi

# *Case Study*

**Analisis faktor paling berpengaruh dalam penjualan Smartphone**

Statistika

Probabilitas

| Variabel | Koefisien | Nilai p | Interpretasi Pengaruh |
|----------|-----------|---------|----------------------|
| Resolusi Kamera | 1.5 | 0.03 | Signifikan |
| RAM | 0.8 | 0.10 | Tidak Signifikan |
| Harga | -2.2 | 0.01 | Signifikan |
| Rating Pelanggan | 2.5 | 0.002 | Sangat Signifikan |
| Jumlah Promosi | 3.0 | 0.04 | Signifikan |

# Case Study

**Analisis faktor paling berpengaruh dalam penjualan Smartphone**

K

**Kesimpulan**:

o Rating pelanggan paling berpengaruh terhadap penjualan smartphone, artinya pelanggan lebih cenderung membeli produk yang memiliki rating tinggi dan ulasan positif.

o Harga yang lebih rendah dan jumlah promosi yang banyak dapat meningkatkan penjualan.

o Resolusi kamera memiliki pengaruh lebih besar dibandingkan RAM.

# Case Study

**Analisis faktor paling berpengaruh dalam penjualan Smartphone**

<div style="border:1px solid #e08a3c; background:#fce4cf; display:inline-block;">
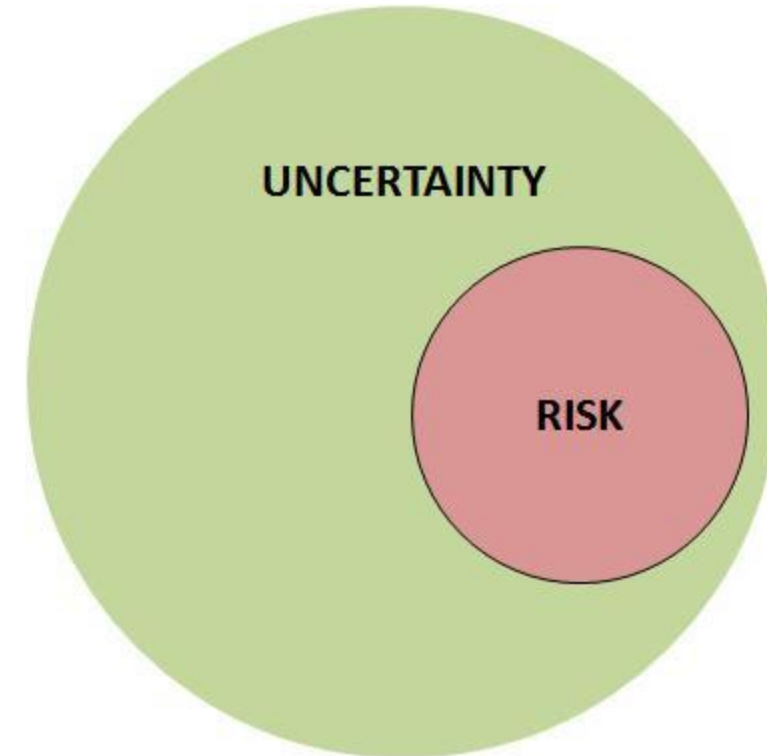
**W**

</div>

**Keputusan**:
o Meningkatkan ulasan pelanggan dan promosi untuk meningkatkan penjualan, misalnya dengan menggunakan jasa para *influencer* sosial media agar meng-*endorse* produk kita.

o Menetapkan harga yang kompetitif untuk mendorong lebih banyak penjualan.

o Memfokuskan promosi pada fitur yang lebih dihargai oleh pelanggan seperti resolusi kamera, dibandingkan RAM.
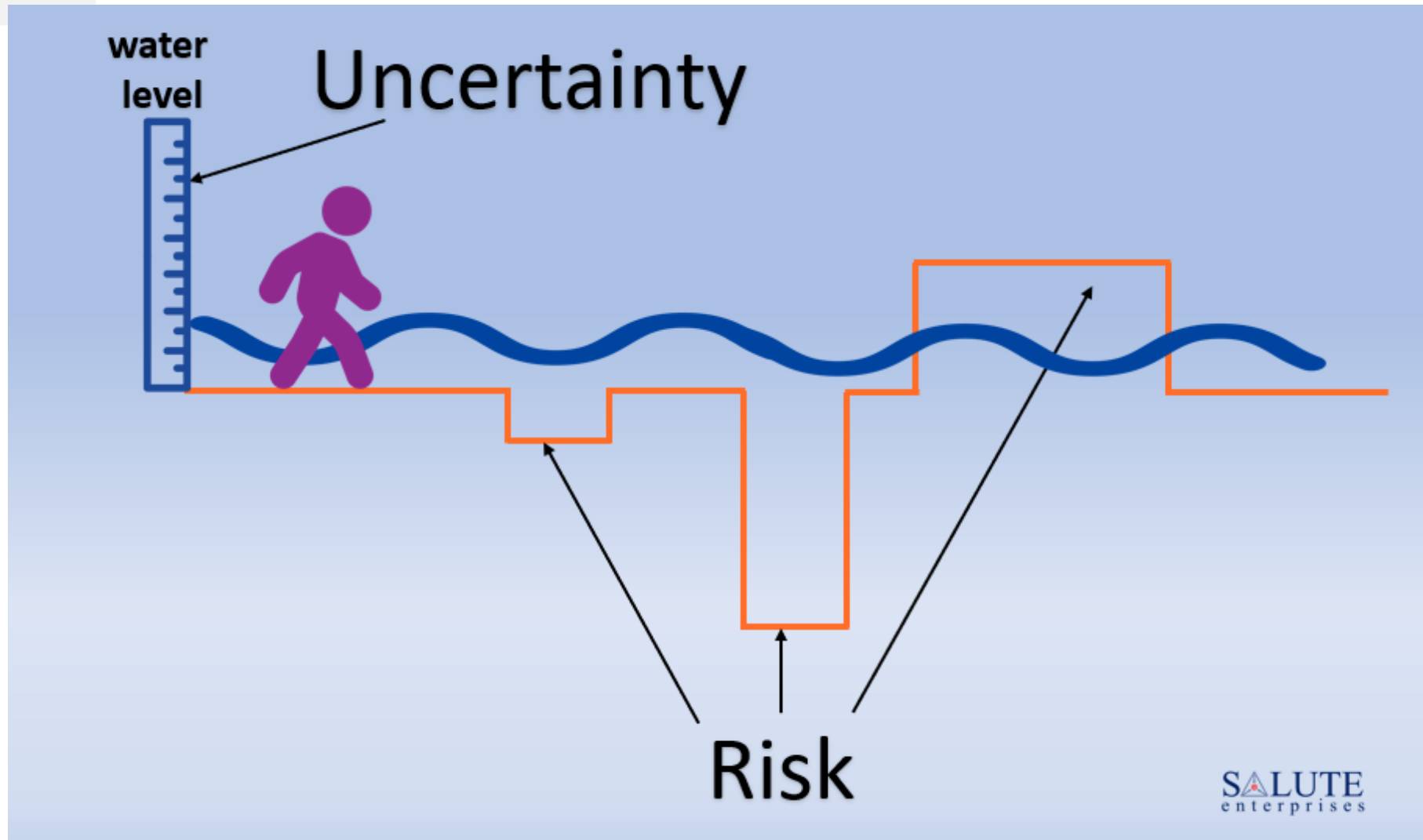
# *Why learn Statistics?*

o **Numerical information is everywhere** and we dealing with **uncertainty**

o Statistical techniques are used to **make decisions** that affect our daily lives

o The knowledge of statistical methods will help you understand how decisions are made and give you a **better understanding** of how they affect you



This Photo by Unknown Author is licensed under CC BY-NC-ND

*INSTITUT TEKNOLOGI SEPULUH NOPEMBER, Surabaya - Indonesia*

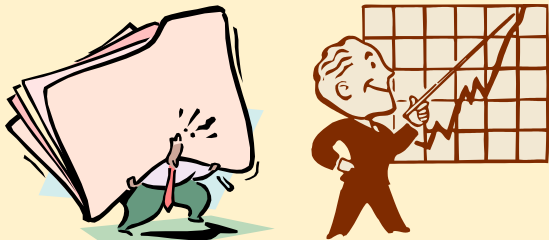# *Why learn Statistics?*

# *Statistics in Information Technology*

o **Data Mining** is the analysis of information in a database, using tools that look for trends or irregularities in large data sets.

o **Data Compression** is the coding of data using compact formulas, called algorithms, and utilities to save storage space or transmission time.

o **Speech Recognition** is the identification of spoken words by a machine. The spoken words are turned into a sequence of numbers and matched against coded dictionaries.

o **Vision and Image Analyses** use statistics to solve contemporary and practical problems in computer vision, image processing, and artificial intelligence.

o **Human/Computer Interaction** uses statistics to design, implement, and evaluate new technologies that are useable, useful, and appealing to a broad cross-section of people.

o **Network/Traffic Modeling** uses statistics to avoid network congestion while fully exploiting the available bandwidth.

o **Stochastic Optimization** uses chance and probability models to develop the most efficient code for finding the solution to a problem.

# The Field of Statistics

# Descriptive Statistics

Descriptive statistics consist of procedures for:

**1**

**Tabulating or graphing** the general characteristics of a set of data.



**2**

**Describing some characteristics of this set** such as measures of central tendency or measures of dispersion.

# Inferential Statistics

Inferential statistics consists of a set of procedures that helps us make **inferences** and **predictions** about a whole population based on information from a **sample of the population**.

o **Estimation / Prediction / Forecasting**
Ex: Estimate the population mean weight using the sample mean weight

o **Hypothesis testing**
Ex: Test the claim that the population means the weight is 120 pounds

o **Make decisions**



Population

Sample

Inferential statistics applied on a sample of the population

The appropriate conclusions regarding the population's features

# Sample and Population

# Sample and Population

○ A **population** consists of the set of **all measurements** for which the investigator is interested.

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

○ A **census** is a **complete enumeration** of every item in a population.



Population (N)

Sample ($n$)

*INSTITUT TEKNOLOGI SEPULUH NOPEMBER, Surabaya - Indonesia*

# *Population Parameter and Sample Statistic*



We want to know about these ...

... but we only have those limited data

Random Selection

**Sample**

**Population**

o A Population **Parameter** is a numerical value that describes a characteristic of an entire population.

o A Sample **Statistic** is a numerical value that describes a characteristic of a sample, which is a subset of the population.

**Parameter**  μ          Inference          x̄  **Statistic**

(Population mean)          (Sample mean)

# *Random Sample*

- Sampling from the population is often done **randomly**, such that every possible sample of equal size (n) will have an equal chance of being selected.

- A sample selected in this way is called a simple random sample or just a **random sample.**

- A random sample allows the chance to determine its elements.



Diastolic Blood Pressure?

Mean = 78 mm Hg

Samples

Mean = 75

Mean = 67

Mean = 71.3

# Why Sample?

o **Census** of a population may be:

    o Impossible

    o Impractical

    o Too costly

# Summary Measures

# Summary Measures

o Statistical values that **summarize** or **describe key characteristics** of a dataset.

o Provide a quick, overall understanding of the data by focusing on aspects such as central tendency, spread, and distribution.

**Describing Data Numerically**

**Central Tendency**
- Arithmetic Mean
- Median
- Mode
- Geometric Mean

**Variation / Dispersion (Spread)**
- Range
- Interquartile Range
- Variance
- Standard Deviation
- Coefficient of Variation

**Shape**
- Skewness
- Kurtosis

**Position**
- Quartiles
- Percentiles

# Measures of Central Tendency or Location

o Statistical tools are used to determine the **center** or **typical value** in a dataset.

o Summarize the data by identifying a **central point** that best represents the **distribution** of values.

| Median | o Middle value when sorted in order of magnitude<br>o 50th percentile |
|--------|-----------------------------------------------------------------------|
| Mode   | o Most frequently-occurring value                                     |
| Mean   | o Average                                                             |

# Measures of Central Tendency or Location

o Statistical tools are used to determine the **center** or **typical value** in a dataset.

o Summarize the data by identifying a **central point** that best represents the **distribution** of values.



(a) Negatively skewed — Mode, Median, Mean; Frequency; Negative direction

(b) Normal (no skew) — Mean, Median, Mode; The normal curve represents a perfectly symmetrical distribution

(c) Positively skewed — Mode, Median, Mean; Positive direction

# Example – Median

| Sales | Sorted Sales |
|-------|--------------|
| 9 | 6 |
| 6 | 9 |
| 12 | 10 |
| 10 | 12 |
| 13 | 13 |
| 15 | 14 |
| 16 | 14 |
| 14 | 15 |
| 14 | 16 |
| 16 | 16 | ← Median |
| 17 | 16 |
| 16 | 17 |
| 24 | 17 |
| 21 | 18 |
| 22 | 18 |
| 18 | 19 |
| 19 | 20 |
| 18 | 21 |
| 20 | 22 |
| 17 | 24 |

## Median
## 50th Percentile

$(20+1)50/100 = 10.5$

**16 + (.5)(0) = 16**

o The **median** is the middle value of data sorted in order of magnitude. It is the 50th percentile.

o Useful for **skewed distributions or datasets with outliers**, as it isn't affected by extreme values.

# *Example – Median*

**Median**

Arrange the observations in ascending order.

**Number of observations ($n$) is odd.**

The median is the middle value, which is at position

$$\left(\frac{n+1}{2}\right)$$

**Number of observations ($n$) is even.**

The median is the average of the two middle values.

1. Find the value at position $\left(\frac{n}{2}\right)$

2. Find the value at position $\left(\frac{n}{2}\right)+1$

3. Find the average of the two values to get the median.

# *Example – Mode*

o The **mode** is the most frequently occurring value. It is the value with the **highest frequency**.

o Best for **categorical data** or when identifying the most common occurrence is important.



**Mode = 16**

# Arithmetic Mean or Average

o The **mean** of a set of observations is their average - the sum of the observed values divided by the number of observations.

o Best for datasets with **evenly distributed data without extreme outliers**

| Population Mean | Sample Mean |
|---|---|
| $$\mu = \frac{\sum_{i=1}^{N} x_i}{N}$$ | $$\overline{X} = \frac{\sum_{i=1}^{n} x_i}{n}$$ |
| $N$ = number of items in the population | $n$ = number of items in the sample |

# *Example – Mean*

| Sales |
|:-----:|
| 9 |
| 6 |
| 12 |
| 10 |
| 13 |
| 15 |
| 16 |
| 14 |
| 14 |
| 16 |
| 17 |
| 16 |
| 24 |
| 21 |
| 22 |
| 18 |
| 19 |
| 18 |
| 20 |
| 17 |
| 317 |

$$\bar{x} = \frac{\sum_{i=1} x}{n} = \frac{317}{20} = 15.85$$

# Measures of Variability or Dispersion (Spread)

- Quantify how much the **data values in a dataset differ from the central value** (mean or median).

- These measures help **describe the distribution's spread** and indicate how concentrated or scattered the data is.

| Range | o Difference between maximum and minimum values |
|---|---|
| **Interquartile Range** | o Difference between third and first quartile (Q3 - Q1) |
| **Variance** | o Average of the squared deviations from the mean <br> o Definitions of population variance and sample variance differ slightly |
| **Standard Deviation** | o Square root of the variance |

# Example – Range and Interquartile Range

| Sales | Sorted Sales | Rank | |
|-------|--------------|------|---|
| 9 | 6 | 1 | ← Minimum |
| 6 | 9 | 2 | |
| 12 | 10 | 3 | |
| 10 | 12 | 4 | |
| 13 | 13 | 5 | |
| 15 | 14 | 6 | ← First Quartile |
| 16 | 14 | 7 | |
| 14 | 15 | 8 | |
| 14 | 16 | 9 | |
| 16 | 16 | 10 | |
| 17 | 16 | 11 | |
| 16 | 17 | 12 | |
| 24 | 17 | 13 | |
| 21 | 18 | 14 | |
| 22 | 18 | 15 | |
| 18 | 19 | 16 | ← Third Quartile |
| 19 | 20 | 17 | |
| 18 | 21 | 18 | |
| 20 | 22 | 19 | |
| 17 | 24 | 20 | ← Maximum |

**Range:** Maximum - Minimum =
24 - 6 = 18

$Q_1 = 13 + (.25)(1) = 13.25$

$Q_3 = 18 + (.75)(1) = 18.75$

**Interquartile Range:** $Q_3 - Q_1 =$
18.75 - 13.25 = 5.5

# *Variance and Standard Deviation*



**population**

11, 2, 3, 5, 8, 4, 9, 12, 7, 6, 10

sample variance:

$$s^2 = \frac{\sum(x_i - \bar{x})^2}{n-1} = 10$$

**population**

2, 4, 6, 8, 10

population variance:

$$\sigma^2 = \frac{\sum(x_i - \bar{x})^2}{n} = 8$$

| **Population Variance** | **Sample Variance** |
|---|---|
| $$\sigma^2 = \frac{\sum_{i=1}^{N}(x_i - \mu)^2}{N}$$ | $$s^2 = \frac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{n-1}$$ |
| $\sigma^2$ = population variance | $s^2$ = sample variance |
| $x_i$ = value of $i^{th}$ element | $x_i$ = value of $i^{th}$ element |
| $\mu$ = population mean | $\bar{x}$ = sample mean |
| $N$ = population size | $n$ = sample size |

# Calculation of Sample Variance

| $x$ | $x - \bar{x}$ | $(x - \bar{x})^2$ | | $x^2$ |
|-----|-----|-----|-----|-----|
| 6 | -9.85 | 97.0225 | 36 | |
| 9 | -6.85 | 46.9225 | 81 | |
| 10 | -5.85 | 34.2225 | 100 | |
| 12 | -3.85 | 14.8225 | 144 | |
| 13 | -2.85 | 8.1225 | | 169 |
| 14 | -1.85 | 3.4225 | | 196 |
| 14 | -1.85 | 3.4225 | 196 | |
| 15 | -0.85 | 0.7225 | 225 | |
| 16 | 0.15 | 0.0225 | 256 | |
| 16 | 0.15 | 0.0225 | 256 | |
| 16 | 0.15 | 0.0225 | 256 | |
| 17 | 1.15 | 1.3225 | 289 | |
| 17 | 1.15 | 1.3225 | 289 | |
| 18 | 2.15 | 4.6225 | 324 | |
| 18 | 2.15 | 4.6225 | 324 | |
| 19 | 3.15 | 9.9225 | 361 | |
| 20 | 4.15 | 17.2225 | 400 | |
| 21 | 5.15 | 26.5225 | 441 | |
| 22 | 6.15 | 37.8225 | 484 | |
| 24 | 8.15 | 66.4225 | 576 | |
| 317 | 0 | 378.5500 | | 5403 |

$$s^2 = \frac{\sum_{i=1}^{n}(x-\bar{x})^2}{(n-1)} = \frac{378.55}{(20-1)}$$

$$= \frac{378.55}{19} = 19.923684$$

$$= \frac{\sum_{i=1}^{n} x^2 - \frac{\left(\sum_{i=1}^{n} x\right)^2}{n}}{(n-1)}$$

$$= \frac{5403 - \frac{317^2}{20}}{(20-1)} = \frac{5403 - \frac{100489}{20}}{19}$$

$$= \frac{5403 - 5024.45}{19} = \frac{378.55}{19} = 19.923684$$

$$s = \sqrt{s^2} = \sqrt{19.923684} = 4.46$$

# Relations between the Mean and Standard Deviation

- The **mean** represents the average value of a dataset or the central point around which the data is distributed.

- The **standard deviation** measures the spread or dispersion of the data points from the mean. It quantifies how far, on average, the data points are from the mean.

- A **low standard deviation** means that the data points are close to the mean (the data is less spread out).

- A **high standard deviation** means that the data points are more spread out from the mean (the data is more dispersed).
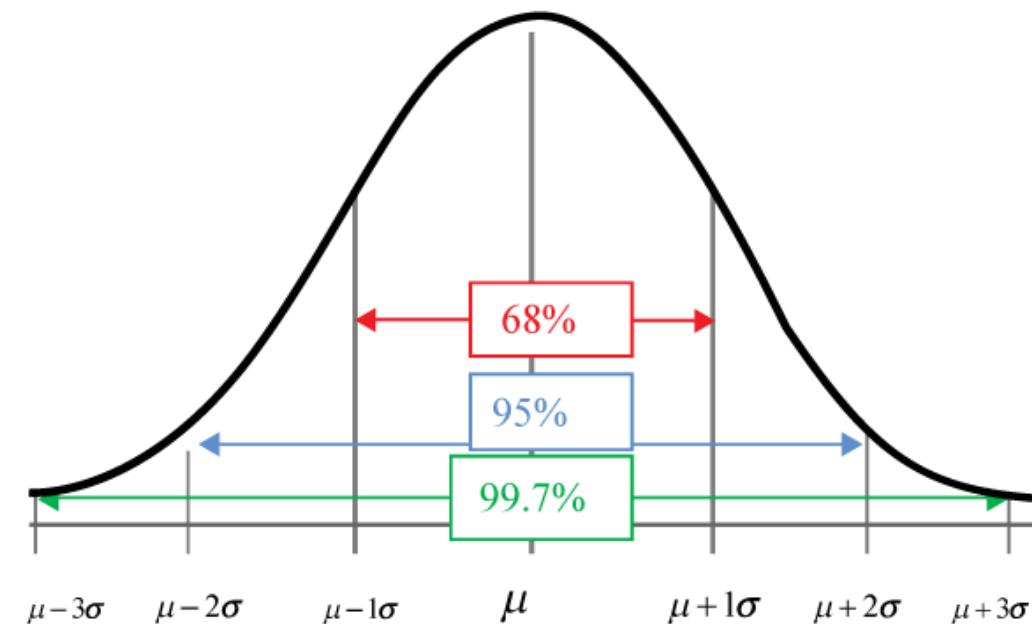
# Relations between the Mean and Standard Deviation

- **Empirical Rule**
  - In a normal distribution (bell-shaped curve), about **68%** of the data points lie within one standard deviation of the mean, **95%** within two standard deviations, and **99.7%** within three standard deviations.
  - This is known as the $68 - 95 - 99.7$ rule or the empirical rule.

**Empirical Rule**
(Normal Distributions)

68%

95%

99.7%

$\mu - 3\sigma \quad \mu - 2\sigma \quad \mu - 1\sigma \quad \mu \quad \mu + 1\sigma \quad \mu + 2\sigma \quad \mu + 3\sigma$
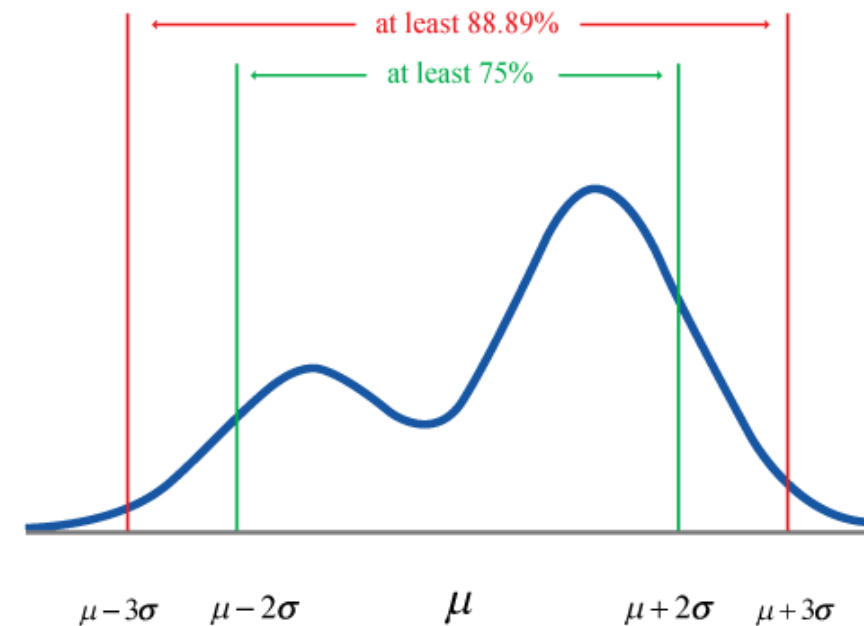
# Relations between the Mean and Standard Deviation

**Assignments**

Pelajari tentang Chebyshev's Inequality ☺



Chebyshev's Inequality
(Any Distribution)

at least 88.89%
at least 75%

$\mu-3\sigma$    $\mu-2\sigma$    $\mu$    $\mu+2\sigma$    $\mu+3\sigma$

# Exercise 1

o Jelaskan dan pelajari tentang **Aturan Empiris (Empirical Rules)**

o Kerjakan:
  o Tinggi badan siswa di sebuah sekolah mengikuti distribusi normal dengan rata-rata 160 cm dan standar deviasi 7 cm. Gunakan Aturan Empiris untuk menjawab pertanyaan berikut:
    o Berapa rentang tinggi badan di mana sekitar 68% siswa berada?
    o Berapa rentang tinggi badan di mana sekitar 95% siswa berada?
    o Berapa rentang tinggi badan di mana sekitar 99.7% siswa berada?

# Exercise 2

o Jelaskan dan pelajari tentang **Aturan Empiris (Empirical Rules)**

o Kerjakan:
   o Dalam sebuah uji coba, waktu reaksi dari sejumlah pengemudi diukur. Diketahui bahwa waktu reaksi rata-rata adalah 0,8 detik dengan standar deviasi 0,1 detik. Berdasarkan Aturan Empiris, tentukan rentang waktu reaksi di mana:
      o 68% pengemudi berada
      o 95% pengemudi berada
      o 99.7% pengemudi berada

# *Exercise 3*

o Jelaskan dan pelajari tentang **Teorema Chebyshev (Chebyshev's Theorem)**

o Kerjakan:
  o Sebuah perusahaan mencatat waktu produksi barang dengan rata-rata 40 menit dan standar deviasi 5 menit. Gunakan Teorema Chebyshev untuk menjawab pertanyaan berikut:
    o Berapa proporsi minimum waktu produksi yang berada dalam jarak 3 standar deviasi dari rata-rata?

# *Exercise 4*

o Jelaskan dan pelajari tentang **Teorema Chebyshev (Chebyshev's Theorem)**

o Kerjakan:
  o Dari sebuah penelitian, diketahui bahwa penghasilan bulanan dari 100 orang karyawan memiliki rata-rata Rp5.000.000 dengan standar deviasi Rp500.000. Tentukan proporsi minimum dari karyawan yang penghasilannya berada dalam jarak 2 standar deviasi dari rata-rata menurut **Teorema Chebyshev**.

# Notation

| Dasar perbandingan | Parameter | Statistik |
|---|---|---|
| Mean | $\mu$ | $\bar{x}$ |
| Standard deviation | $\sigma$ | $s$ |
| Proporsi | $P$ | $\hat{p}$ |
| Elemen data | $X$ | $x$ |
| Ukuran sampel | $N$ | $n$ |
| Koefisiensi korelasi | $\rho$ | $r$ |

*Sumber: Key Differences*

*https://revou.co/revoupedia/kosakata*

# Notation

| Parameter name | Population parameter symbol | Sample statistic |
| --- | --- | --- |
| Number of cases | N | n |
| Mean | $\mu$ (mu) | $\bar{x}$ (Sample mean) |
| Proportion | $\pi$ (Pi) | P (Sample proportion) |
| Variance | $\sigma^2$ (Sigma-square) | $s^2$ (Sample variance) |
| Standard deviation | $\sigma$ (Sigma) | s (sample standard deviation) |
| Correlation | $\rho$ (rho) | r (Sample correlation) |
| Regression Coefficient | $\beta$ (beta) | b (sample regression coefficient) |