

# INTRODUCTION TO PROBABILITY AND STATISTICS

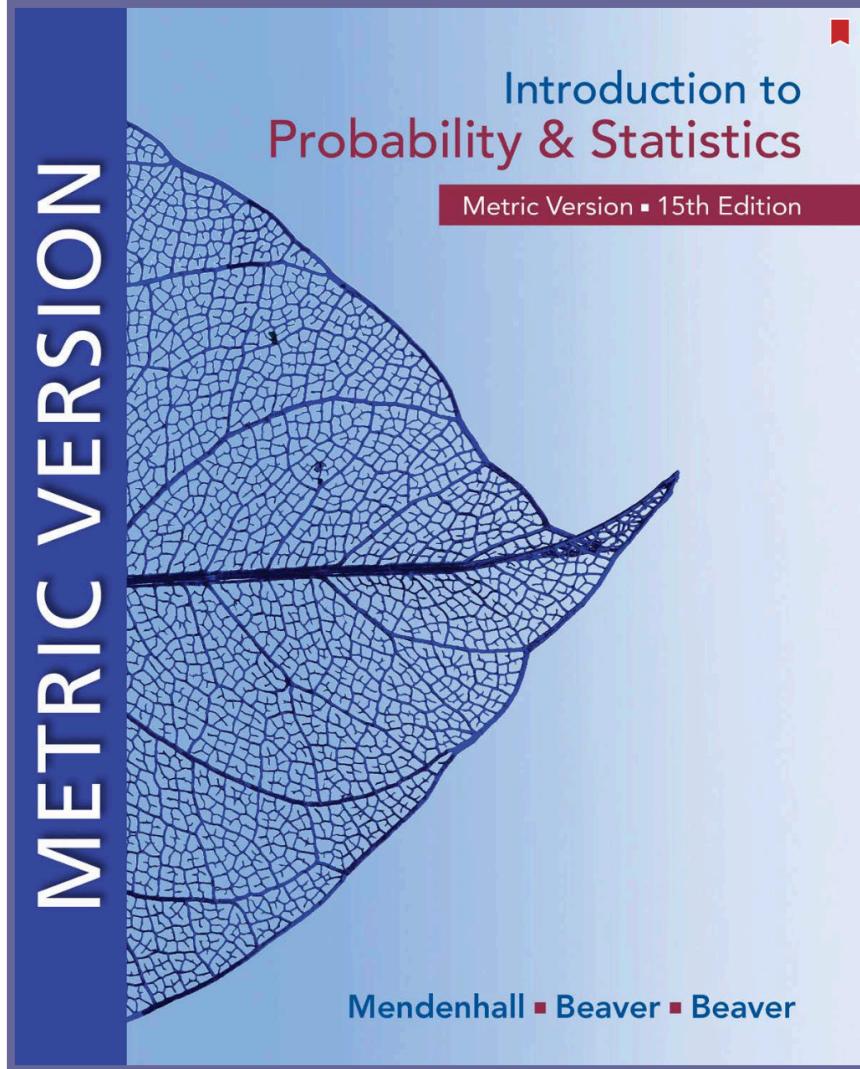
## FOURTEENTH EDITION

William Mendenhall, III • Robert J. Beaver •  
Barbara M. Beaver



1

# METRIC VERSION

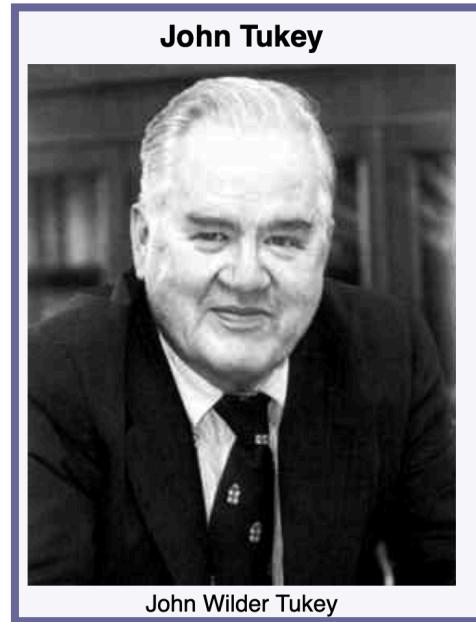


# TEXTBOOK



# EXPLORATORY DATA ANALYSIS

- Tukey (1977, preface): “It is important to understand what you CAN DO before you learn to measure how WELL you seem to have DONE it... Exploratory data analysis can never be the whole story, but nothing else can serve as the foundation stone –as **the first step.**”



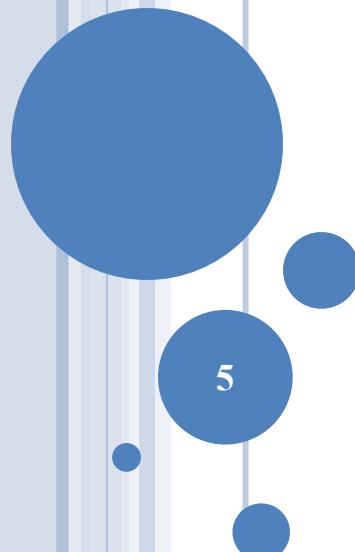
1915-2000

# MORE ON EDA

- ChatGPT:
  - EDA is an initial step in data analysis to summarize main characteristics, often using **visual** methods.
  - It helps in understanding **data structure**, detecting **outliers**, and **identifying patterns, relationships, and anomalies**.
  - EDA provides **insights** for more sophisticated analyses and guides **hypothesis formulation**.
- Two approaches
  - **Numerical measures**: mean, std, max, min, median, etc
  - **Graphs**: pieplot, boxplot, histogram, scatter plots, etc
- **Chapters 1 to 3** are fundamental tools in EDA.

# INTRODUCTION TO PROBABILITY AND STATISTICS

## FOURTEENTH EDITION



### Chapter 1

## Describing Data with Graphs

# VARIABLES AND DATA

- A **variable** is a characteristic that changes or varies over time and/or for different individuals or objects under consideration.
- **Examples:**
  - Hair color: Qualitative, categorical
  - white blood cell count: Quantitative, discrete (integer)
  - time to failure of a computer component: Quantitative, continuous (positive continuous)

# DEFINITIONS

- An **experimental unit** is the individual or object on which a variable is measured.
- A **measurement** results when a variable is actually measured on an experimental unit.
- A set of measurements, called **data**, can be either a **sample** or a **population**.

# Example

- **Variable**
  - Hair color
- **Experimental unit**
  - Person
- **Typical Measurements**
  - Brown, black, blonde, etc.



## EXAMPLE

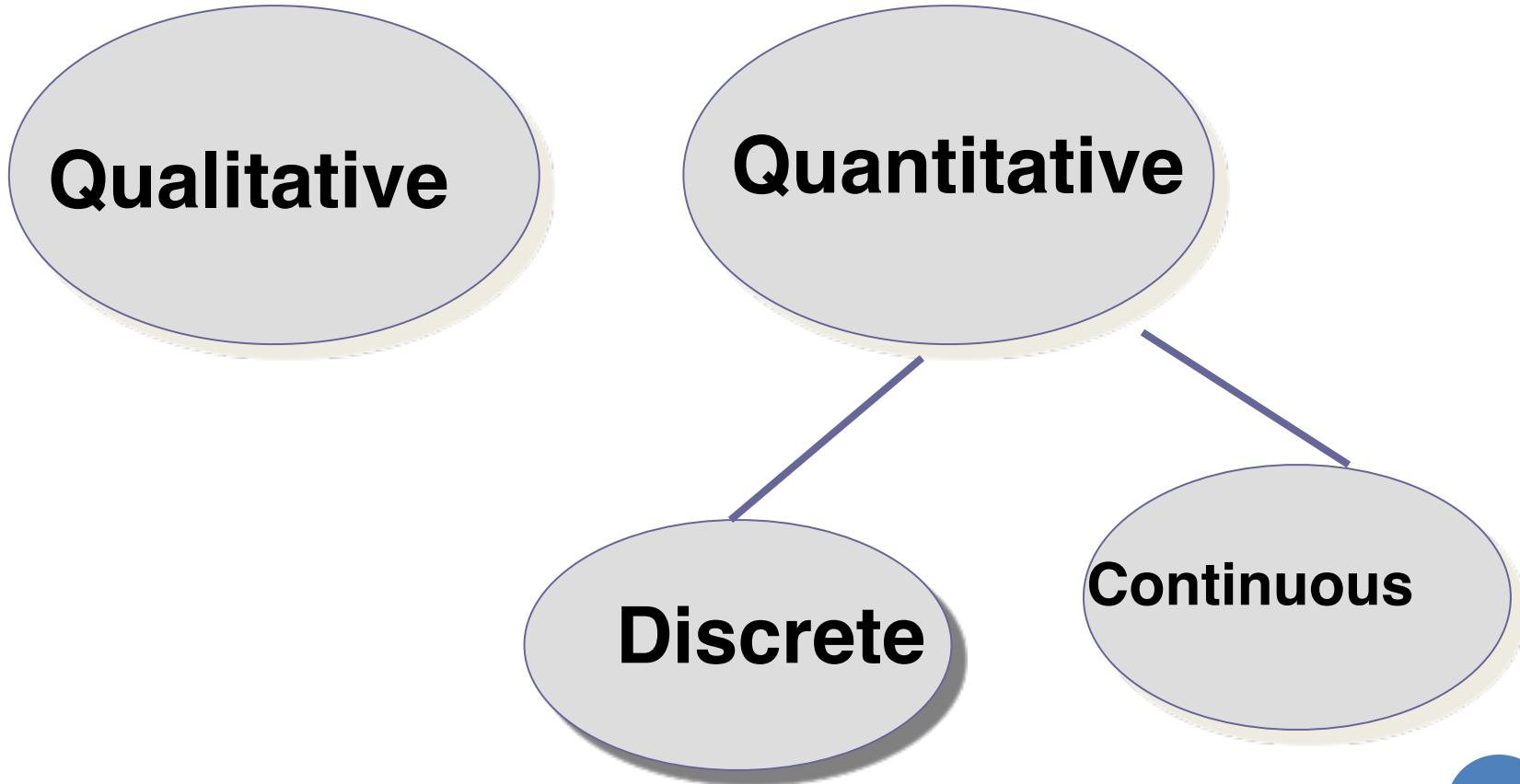
- **Variable**
  - Time until a light bulb burns out
- **Experimental unit**
  - Light bulb
- **Typical Measurements**
  - 1500 hours, 1535.5 hours, etc.



# HOW MANY VARIABLES HAVE YOU MEASURED?

- **Univariate data:** One variable is measured on a single experimental unit.
- **Bivariate data:** Two variables are measured on a single experimental unit.
- **Multivariate data:** More than two variables are measured on a single experimental unit.

# TYPES OF VARIABLES



# TYPES OF VARIABLES

- **Qualitative variables** measure a quality or characteristic on each experimental unit.
- **Examples:**
  - Hair color (black, brown, blonde...)
  - Make of car (Dodge, Honda, Ford...)
  - Gender (male, female)
  - State of birth (California, Arizona,....)

# TYPES OF VARIABLES

• **Quantitative variables** measure a numerical quantity on each experimental unit.

✓ **Discrete** if it can assume only a finite or countable number of values.

- For each orange tree in a grove, the number of oranges is measured.
- For a particular day, the number of cars entering a college campus is measured.



✓ **Continuous** if it can assume the infinitely many values corresponding to the points on a line interval.

- Time until a light bulb burns out

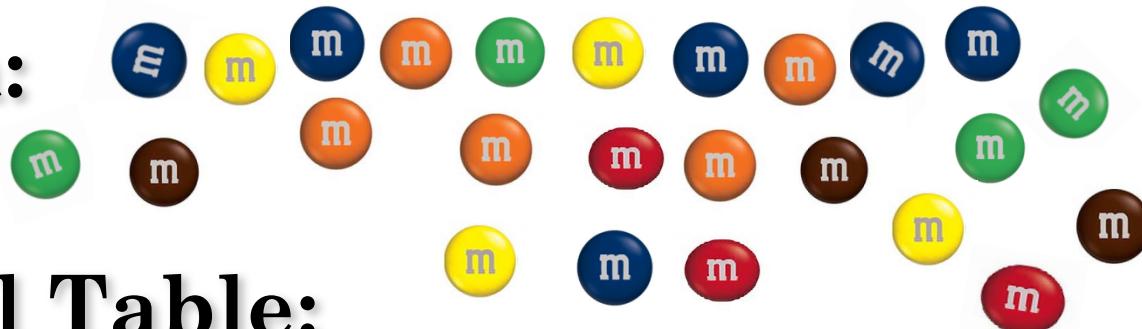
# GRAPHING QUALITATIVE VARIABLES

- Use a **data distribution** to describe:
  - **What values** of the variable have been measured
  - **How often** each value has occurred
- “How often” can be measured 3 ways:
  - Frequency
  - Relative frequency = Frequency/n
  - Percent =  $100 \times$  Relative frequency

# EXAMPLE

- A bag of M&Ms contains 25 candies:

- Raw Data:

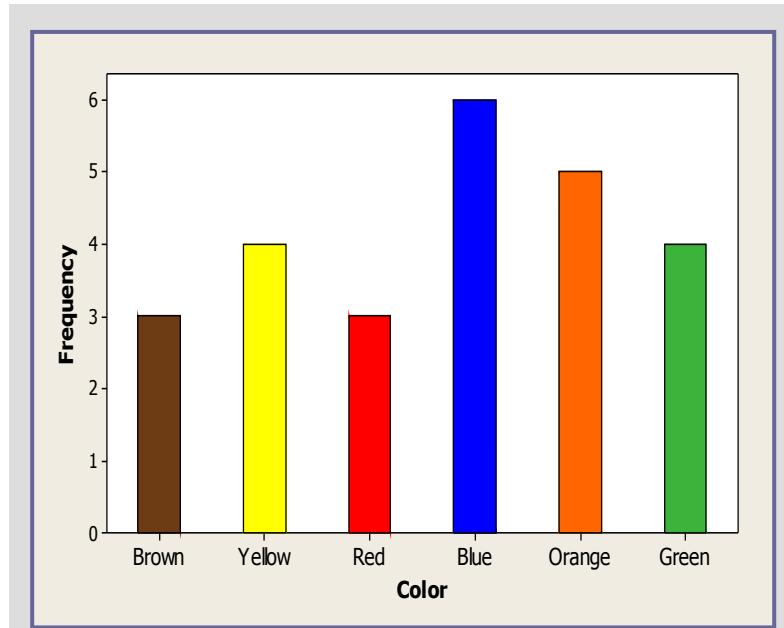


- Statistical Table:

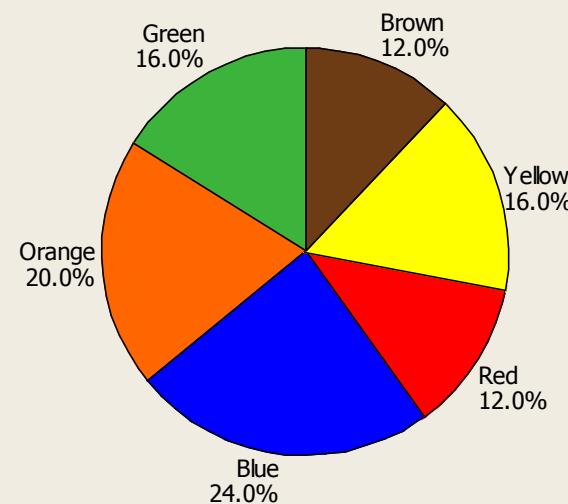
Color	Tally	Frequency	Relative Frequency	Percent
Red		3	$3/25 = .12$	12%
Blue		6	$6/25 = .24$	24%
Green		4	$4/25 = .16$	16%
Orange		5	$5/25 = .20$	20%
Brown		3	$3/25 = .12$	12%
Yellow		4	$4/25 = .16$	16%

# GRAPHS

## Bar Chart



## Pie Chart

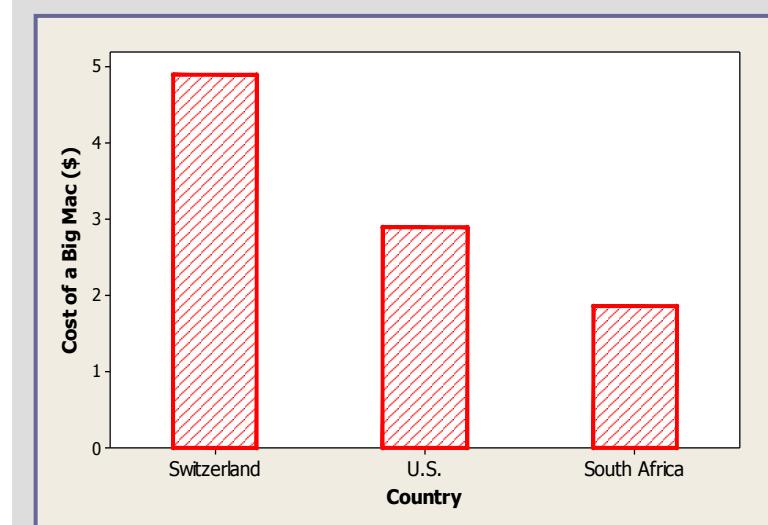


## 2. Quantitative, simple cases

# GRAPHING QUANTITATIVE VARIABLES

- A single quantitative variable measured for different population segments or for different categories of classification can be graphed using a **pie** or **bar** chart.

A Big Mac hamburger costs \$4.90 in Switzerland, \$2.90 in the U.S. and \$1.86 in South Africa.

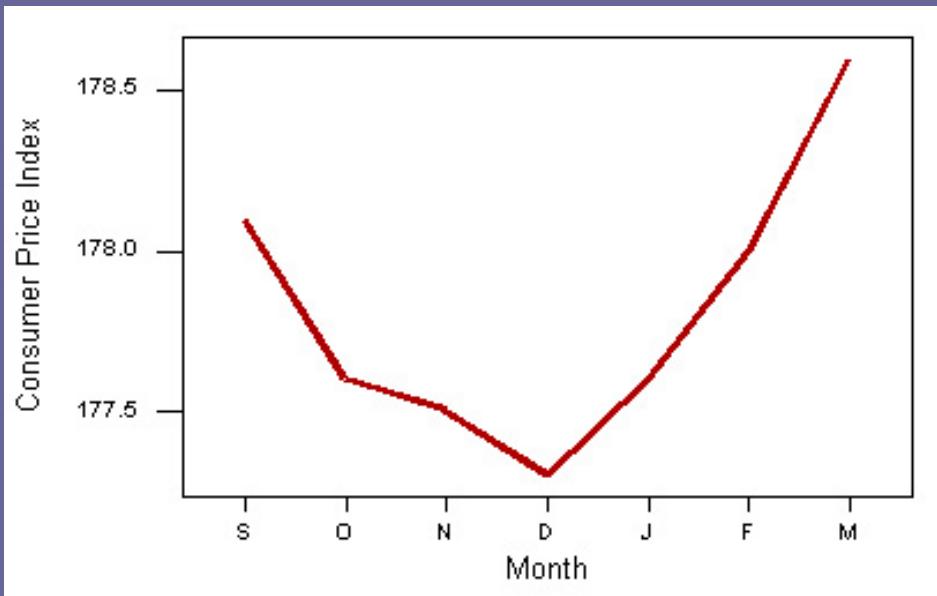


- A single quantitative variable measured over time is called a **time series**. It can be graphed using a **line** or **bar chart**.

### CPI: All Urban Consumers-Seasonally Adjusted

Sept	Oct	Nov	Dec	Jan	Feb	Mar
178.10	177.60	177.50	177.30	177.60	178.00	178.60

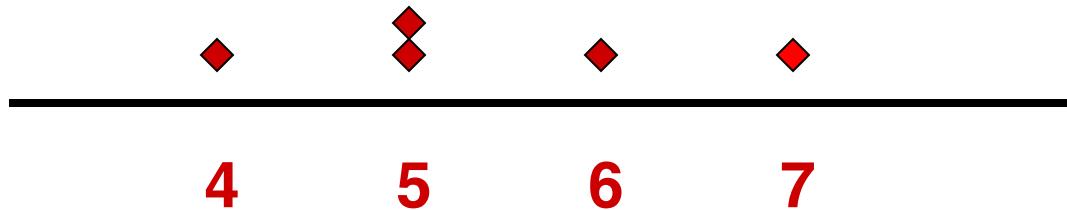
BUREAU OF LABOR STATISTICS



## 2. Quantitative: Discrete

### DOTPLOTS

- The simplest graph for quantitative data
- Plots the measurements as points on a horizontal axis, stacking the points that duplicate existing points.
- **Example:** The set 4, 5, 5, 7, 6



## STEM AND LEAF PLOTS

- A simple graph for quantitative data
- Uses the actual numerical values of each data point.

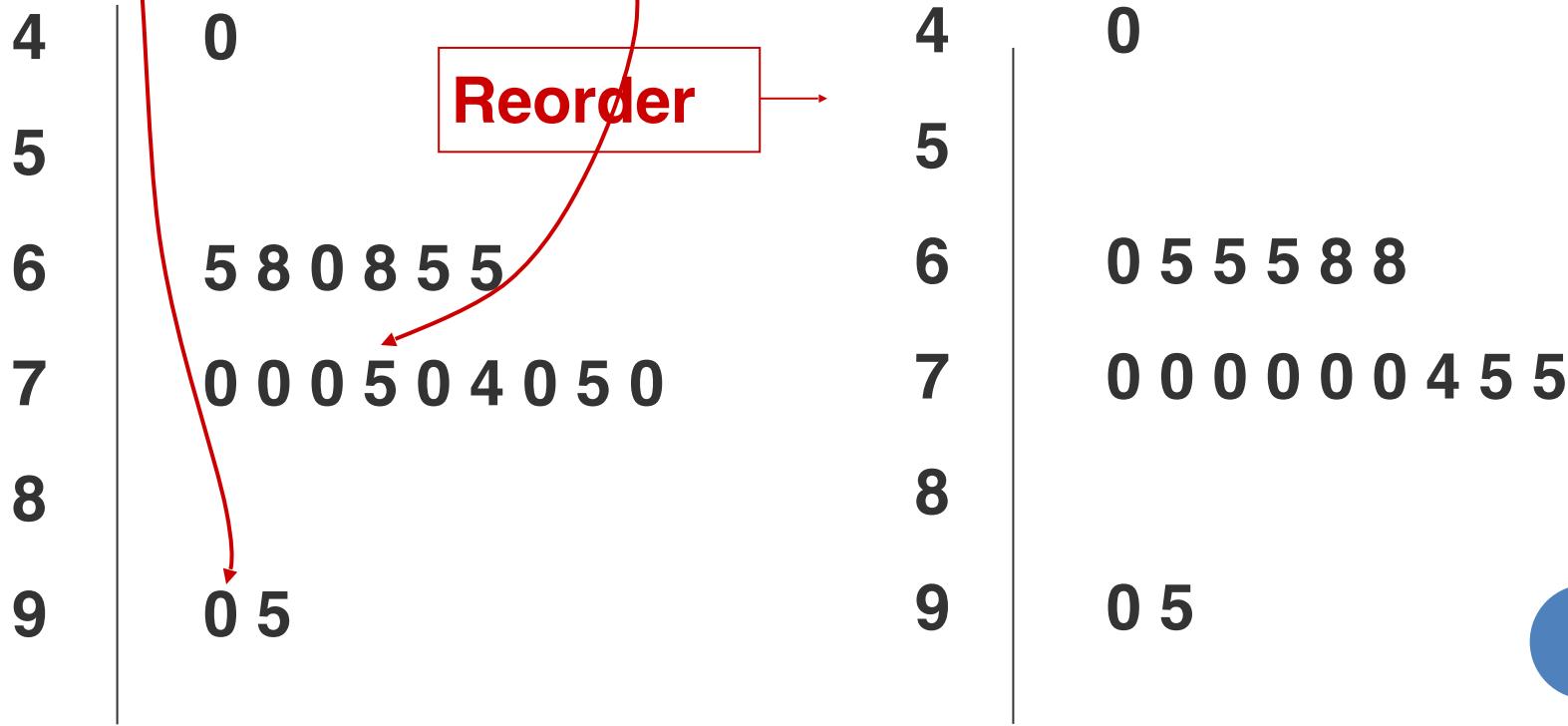
- Divide each measurement into two parts: the **stem** and the **leaf**.
- List the stems in a column, with a **vertical line** to their right.
- For each measurement, record the leaf portion in the **same row** as its matching stem.
- Order** the leaves from lowest to highest in each stem.
- Provide a **key** to your coding.



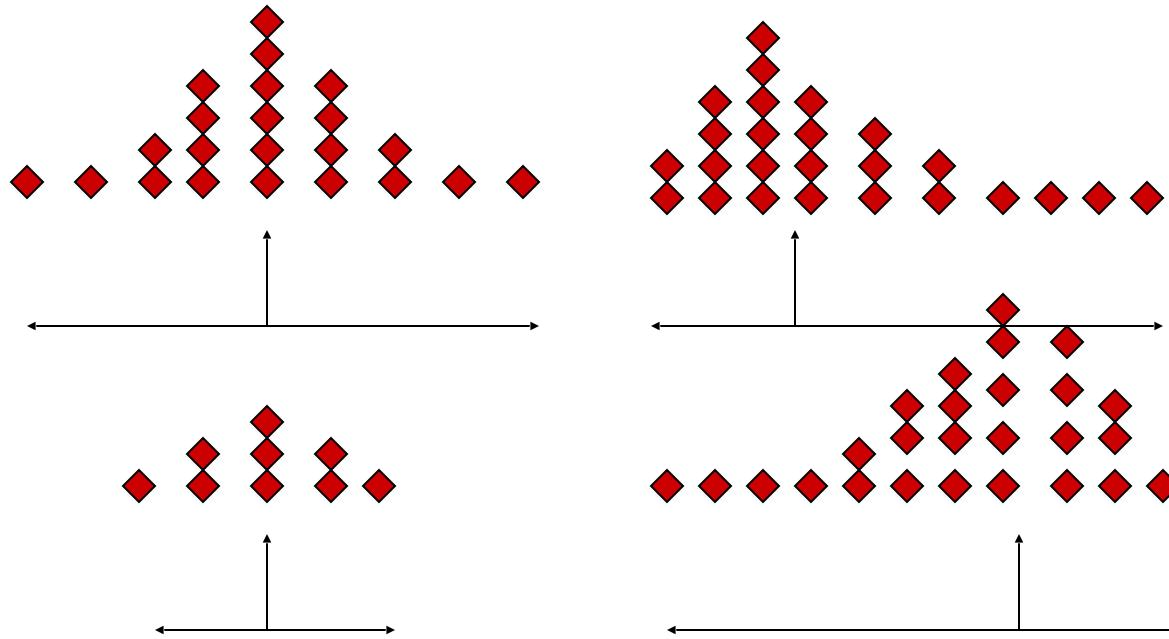
## EXAMPLE

The prices (\$) of 18 brands of walking shoes:

90	70	70	70	75	70	65	68	60
74	70	95	75	70	68	65	40	65

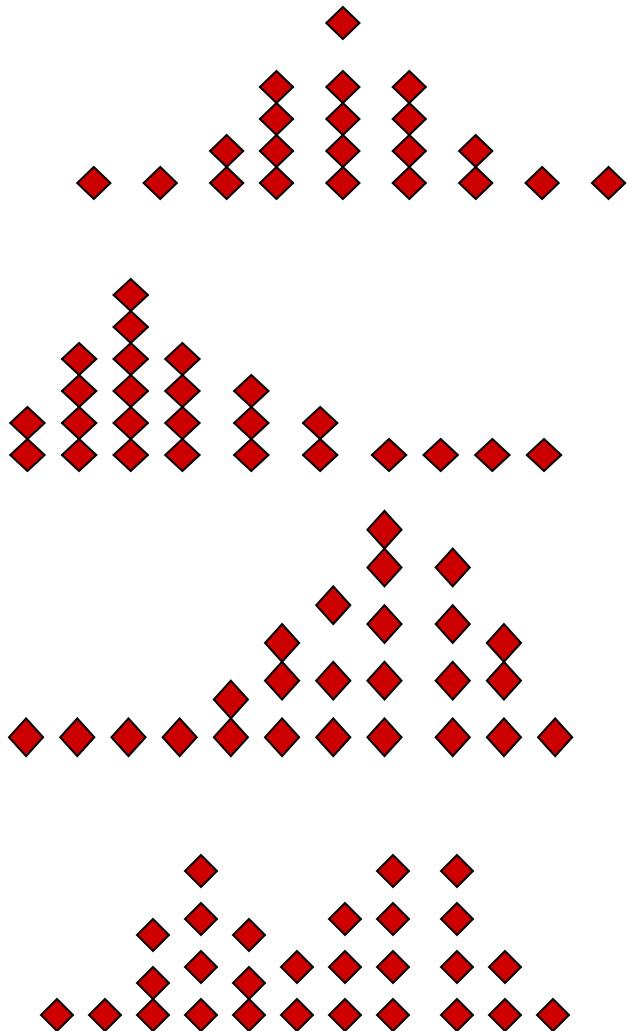


# INTERPRETING GRAPHS: LOCATION AND SPREAD



- Where is the data centered on the horizontal axis, and how does it spread out from the center?

# INTERPRETING GRAPHS: SHAPES



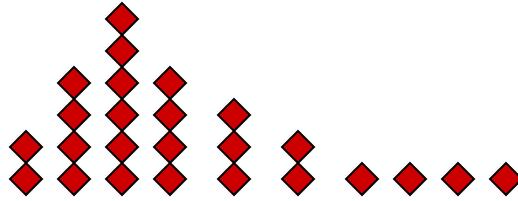
**Mound shaped and symmetric (mirror images)**

**Skewed right: a few unusually large measurements**

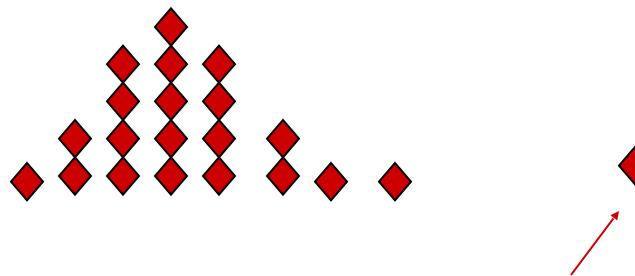
**Skewed left: a few unusually small measurements**

**Bimodal: two local peaks**

# INTERPRETING GRAPHS: OUTLIERS

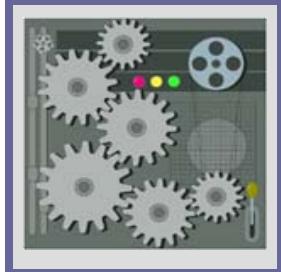


No Outliers



Outlier

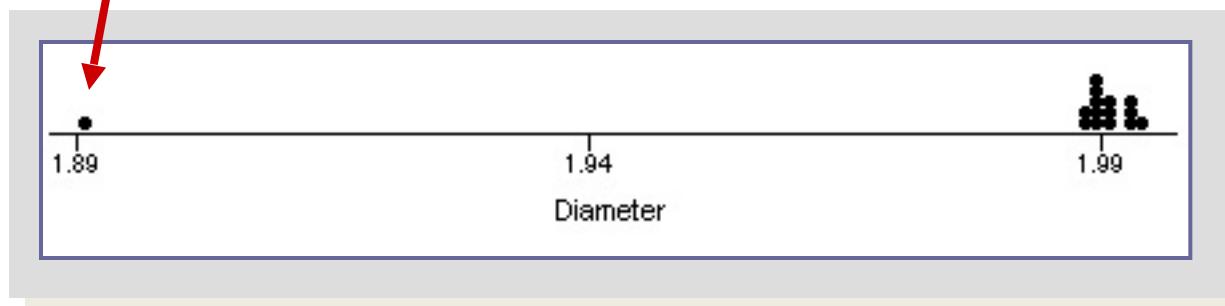
- Are there any strange or unusual measurements that stand out in the data set?



# EXAMPLE

- A quality control process measures the diameter of a gear being made by a machine (cm). The technician records 15 diameters, but inadvertently makes a typing mistake on the second entry.

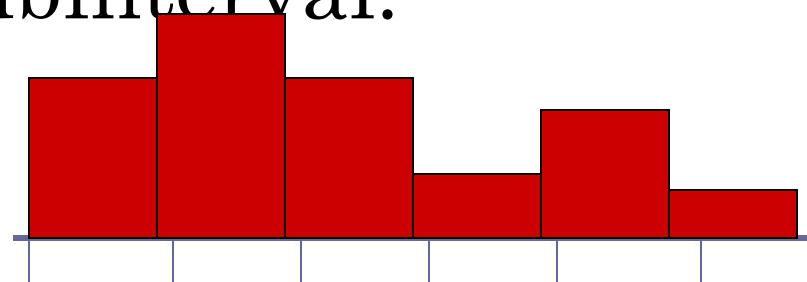
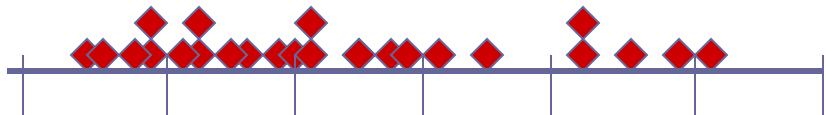
1.991	<b>1.891</b>	1.991	1.988	1.993	1.989	1.990	1.988
1.988	<b>1.993</b>	1.991	1.989	1.989	1.993	1.990	1.994



### 3. Quantitative: Continuous

## RELATIVE FREQUENCY HISTOGRAMS

- A **relative frequency histogram** for a quantitative data set is a bar graph in which the height of the bar shows “how often” (measured as a proportion or relative frequency) measurements fall in a particular class or subinterval.



# RELATIVE FREQUENCY HISTOGRAMS

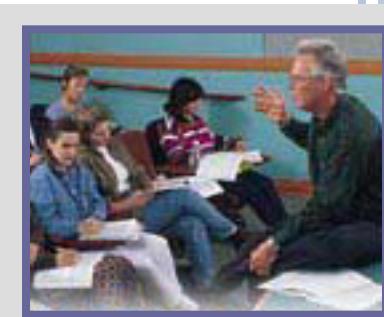
- Divide the range of the data into **5-12 subintervals** of equal length.
- Calculate the **approximate width** of the subinterval as Range/number of subintervals.
- Round the approximate width up to a convenient value.
- Use the method of **left inclusion** including the left endpoint, but not the right in your tally.
- Create a **statistical table** including the subintervals, their frequencies and relative frequencies.

# RELATIVE FREQUENCY HISTOGRAMS

- Draw the **relative frequency histogram** plotting the subintervals on the horizontal axis and the relative frequencies on the vertical axis.
- The height of the bar represents
  - The **proportion** of measurements falling in that class or subinterval.
  - The **probability** that a single measurement, drawn at random from the set, will belong to that class or subinterval.

# EXAMPLE

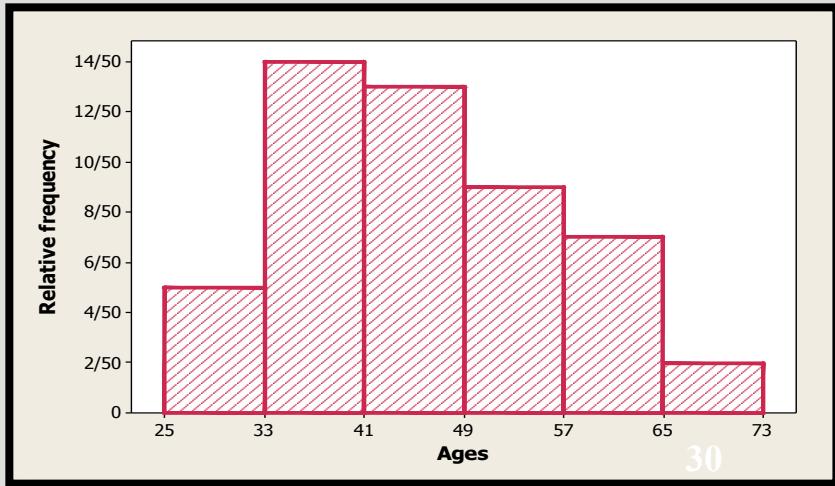
The ages of 50 tenured faculty at a state university.



◦ 34	48	<b>70</b>	63	52	52	35	50	37	43	53	43	52	44
◦ 42	31	36	48	43	<b>26</b>	58	62	49	34	48	53	39	45
◦ 34	59	34	66	40	59	36	41	35	36	62	34	38	28
◦ 43	50	30	43	32	44	58	53						

- We choose to use **6** intervals.
- Minimum class width =  $(70 - 26)/6 = 7.33$
- Convenient class width = **8**
- Use **6** classes of length **8**, starting at **25**.

Age	Tally	Frequency	Relative Frequency	Percent
25 to < 33	<del>1111</del>	5	$5/50 = .10$	10%
33 to < 41	<del>1111 1111 1111</del>	14	$14/50 = .28$	28%
41 to < 49	<del>1111 1111 111</del>	13	$13/50 = .26$	26%
49 to < 57	<del>1111 1111</del>	9	$9/50 = .18$	18%
57 to < 65	1111 11	7	$7/50 = .14$	14%
65 to < 73	11	2	$2/50 = .04$	4%



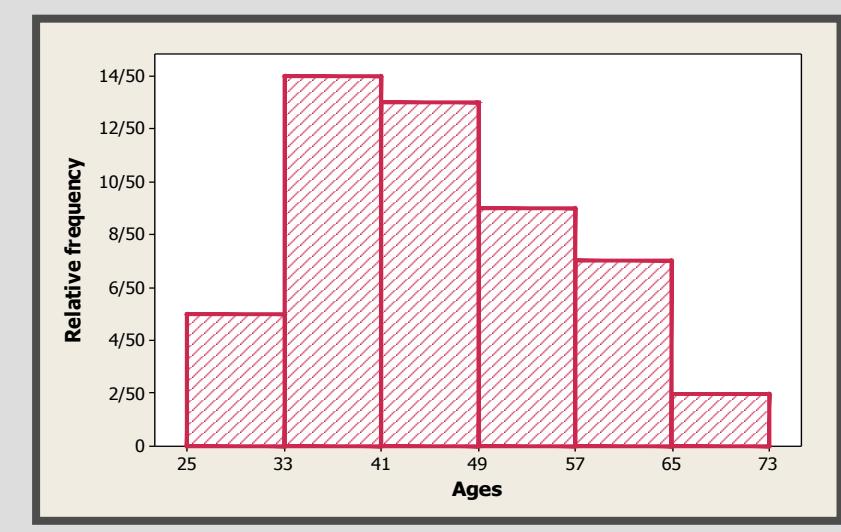
# Describing the Distribution

Shape? **Skewed right**

Outliers? **No.**

What proportion of the tenured faculty are younger than 41?

What is the probability that a randomly selected faculty member is 49 or older?



$$(14 + 5)/50 = 19/50 = .38$$

$$(9 + 7 + 2)/50 = 18/50 = .36$$



# KEY CONCEPTS

## I. How Data Are Generated

1. Experimental units, variables, measurements
2. Samples and populations
3. Univariate, bivariate, and multivariate data

## II. Types of Variables

1. Qualitative or categorical
2. Quantitative
  - a. Discrete
  - b. Continuous

## III. Graphs for Univariate Data Distributions

1. Qualitative or categorical data
  - a. Pie charts
  - b. Bar charts

# KEY CONCEPTS

## 2. Quantitative data

- a. Pie and bar charts
- b. Line charts
- c. Dotplots
- d. Stem and leaf plots
- e. Relative frequency histograms

## 3. Describing data distributions

- a. Shapes—symmetric, skewed left, skewed right, unimodal, bimodal
- b. Proportion of measurements in certain intervals
- c. Outliers