# Descriptive Statistics

(Slides used with permission)

Author: Kristin L. Sainani, PhD
Associate Professor with Health Research and Policy at Stanford University
Webpage: https://web.stanford.edu/~kcobb/

1

# Types of Variables: Overview

**Categorical**          **Quantitative**

**binary**   **nominal**   **ordinal**   **discrete**   **continuous**

2 categories +

more categories +

order matters +

numerical  +

uninterrupted

2

# Categorical Variables

- Also known as "qualitative."
- <u>Dichotomous (binary)</u> –  two levels

  - Dead/alive
  - Treatment/placebo
  - Disease/no disease
  - Exposed/Unexposed
  - Heads/Tails
  - Pulmonary Embolism (yes/no)
  - Male/female

3

## Categorical Variables

- <u>Nominal variables</u> – Named categories Order doesn't matter!

  - The blood type of a patient (O, A, B, AB)
  - Marital status
  - Occupation

4

## Categorical Variables

- <u>Ordinal variable</u> – Ordered categories. Order matters!

  - Staging in breast cancer as I, II, III, or IV
  - Birth order—1st, 2nd, 3rd, etc.
  - Letter grades (A, B, C, D, F)
  - Ratings on a scale from 1-5
  - Ratings on: always; usually; many times; once in a while; almost never; never
  - Age in categories (10-20, 20-30, etc.)
  - Shock index categories (Kline et al.)

5

## Quantitative Variables

- Numerical variables; may be arithmetically manipulated.

  - Counts
  - Time
  - Age
  - Height

6

## Quantitative Variables

- <u>Discrete Numbers</u> – a limited set of distinct values, such as whole numbers.

  - Number of new AIDS cases in CA in a year (counts)
  - Years of school completed
  - The number of children in the family (cannot have a half a child!)
  - The number of deaths in a defined time period (cannot have a partial death!)
  - Roll of a die

7

## Quantitative Variables

- <u>Continuous Variables</u> - Can take on any number within a defined range.

  - Time-to-event (survival time)
  - Age
  - Blood pressure
  - Serum insulin
  - Speed of a car
  - Income
  - Shock index (Kline et al.)

8

## Looking at Data

- ✓ How are the data distributed?
  - Where is the center?
  - What is the range?
  - What's the shape of the distribution (e.g., Gaussian, binomial, exponential, skewed)?

- ✓ Are there "outliers"?

- ✓ Are there data points that don't make sense?

9

The first rule of statistics:
USE COMMON SENSE!

90% of the information is
contained in the graph.

10

# Frequency Plots (univariate)

**Categorical variables**
- Bar Chart

**Continuous variables**
- Box Plot
- Histogram

11

# Bar Chart

- Used for categorical variables to show frequency or proportion in each category.
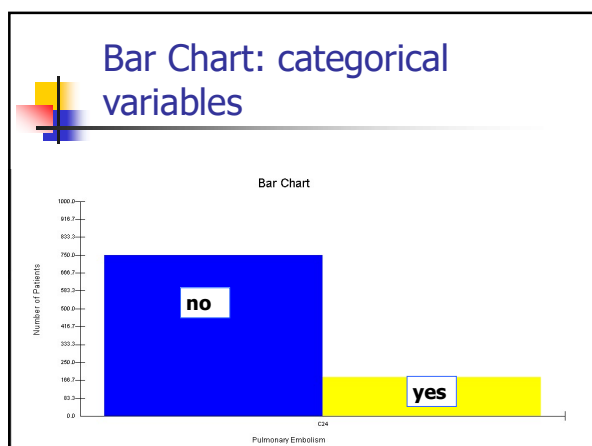- Translate the data from frequency tables into a pictorial representation...

12

## Bar Chart: categorical variables

Bar Chart



**13**

## Bar Chart for SI categories



Much easier to extract information from a bar chart than from a table!

**14**

## Box plot and histograms: for continuous variables

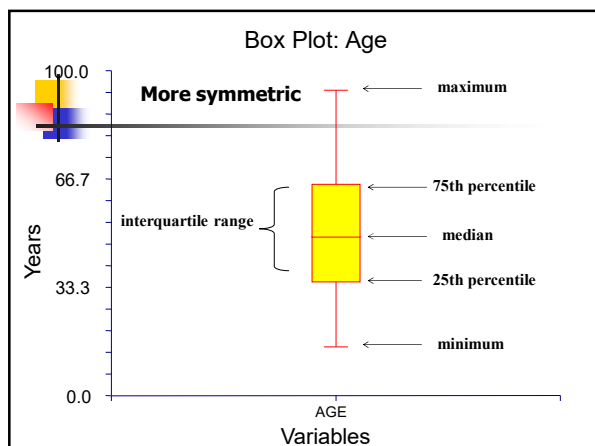- To show the <u>distribution</u> (shape, center, range, variation) of continuous variables.

**15**

5

**16**



**17**



**18**

**19**



**20**



**21**

## Histogram: Age

**Not skewed, but not bell-shaped either…**

Percent

14.0

9.3

4.7

0.0

0.0          33.3          66.7          100.0

AGE (Years)

**22**

# Measures of central tendency

- Mean
- Median
- Mode

**23**

# Central Tendency

- <u>Mean</u> – the average; the balancing point

  *calculation:* the sum of values divided by the sample size

In math shorthand:

$$\overline{X} = \frac{\sum\limits_{i=1}^{n} x}{n} = \frac{X_1 + X_2 + \cdots + X_n}{n}$$

**24**

## Mean: example

Some data:
Age of participants: 17   19   21   22   23   23   23   38

$$\overline{X} = \frac{\sum_{i=1}^{n} X_i}{n} = \frac{17 + 19 + 21 + 22 + 23 + 23 + 38}{8} = 23.25$$

**25**

## Mean of age in some data

Means Section of AGE

| Parameter | Mean | Median | Sum | Mode |
|-----------|------|--------|-----|------|
| Value | 50.19334 | 49 | 46730 | 49 |



**26**

## Mean of age in some data



The balancing point

**27**

## Mean

- The mean is affected by extreme values (outliers)



Mean = 3

$$\frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$

Mean = 4

$$\frac{1+2+3+4+10}{5} = \frac{20}{5} = 4$$

**28**

---

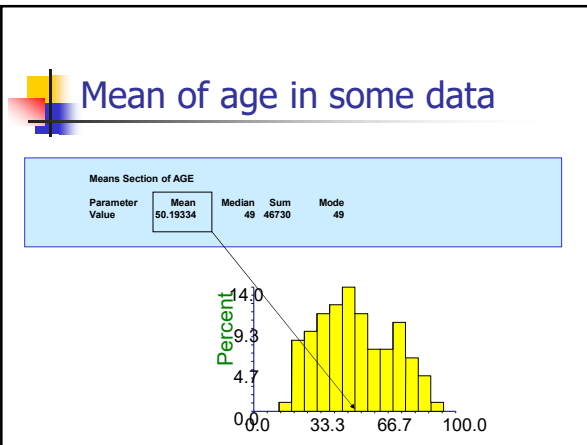## Central Tendency

- <u>Median</u> – the exact middle value

*Calculation:*
- If there are an odd number of observations, find the middle value
- If there are an even number of observations, find the middle two values and average them.

**29**

---

## Median: example

<u>Some data</u>:
Age of participants: 17   19   21   <u>22   23</u>   23   23   38

**Median = (22+23)/2 = 22.5**

**30**

## Median of age in some data

Means Section of AGE

| Parameter | Mean | Median | Sum | Mode |
|---|---|---|---|---|
| Value | 50.19334 | 49 | 46730 | 49 |



31

## Median

- The median is not affected by extreme values (outliers).



Median = 3          Median = 3

▪SSlide from: Statistics for Managers Using Microsoft® Excel 4th Edition, 2004 Prentice-Hall

32

## Central Tendency

- <u>Mode</u> – the value that occurs most frequently

33

## Mode: example

Some data:
Age of participants: 17  19  21  22  23  23  23  38

**Mode = 23  (occurs 3 times)**

## Measures of Variation/Dispersion

- Range
- Percentiles/quartiles
- Interquartile range
- Standard deviation/Variance

## Range

- Difference between the largest and the smallest observations.

Range of age: 94 years-15 years = 79 years



**37**

## Quartiles

| 25% | 25% | 25% | 25% |
|-----|-----|-----|-----|

$Q_1$    $Q_2$    $Q_3$

- The first quartile, $Q_1$, is the value for which 25% of the observations are smaller and 75% are larger
- $Q_2$ is the same as the median (50% are smaller, 50% are larger)
- Only 25% of the observations are greater than the third quartile

**38**

## Interquartile Range

- Interquartile range = 3rd quartile – 1st quartile = $Q_3 - Q_1$

**39**

13

## Interquartile Range: age

| minimum | Q1 | Median (Q2) | Q3 | maximum |
|---------|-----|-------------|-----|---------|
| | 25% | 25% | 25% | 25% |
| 15 | 35 | 49 | 65 | 94 |

Interquartile range
= 65 – 35 = 30

**40**

## Variance

- Average (roughly) of squared deviations of values from the mean

$$S^2 = \frac{\sum_{i}^{n}(x_i - \overline{X})^2}{n-1}$$

**41**

## Why squared deviations?

- Adding deviations will yield a sum of 0.
- Absolute values are tricky!
- Squares eliminate the negatives.

- Result:
  - Increasing contribution to the variance as you go farther from the mean.

**42**

14

## Standard Deviation

- Most commonly used measure of variation
- Shows variation about the mean
- Has the same units as the original data

$$S = \sqrt{\frac{\sum_{i}^{n}(x_i - \overline{X})^2}{n-1}}$$

43

## Calculation Example:
## Sample Standard Deviation

**Age data (n=8) :** 17   19   21   22   23   23   23   38

n = 8        Mean = $\overline{X}$ = 23.25

$$S = \sqrt{\frac{(17-23.25)^2 + (19-23.25)^2 + \cdots + (38-23.25)^2}{8-1}}$$

$$= \sqrt{\frac{280}{7}} = 6.3$$

44

## Comparing Standard Deviations

Data A

11  12  13  14  15  16  17  18  19  20  21

Mean = 15.5
S = 3.338

Data B

11  12  13  14  15  16  17  18  19  20  21

Mean = 15.5
S = 0.926

Data C

11  12  13  14  15  16  17  18  19  20  21

Mean = 15.5
S = 4.570

- SSlide from: Statistics for Managers Using Microsoft® Excel  4th Edition, 2004 Prentice-Hall

45

## Symbol Clarification

- S = <u>Sample</u> standard deviation (example of a "sample statistic")
- $\sigma$ = Standard deviation of the entire population (example of a "population parameter") or from a theoretical probability distribution
- $\bar{X}$ = <u>Sample</u> mean
- $\mu$ = Population or theoretical mean

46

## **The beauty of the normal (bell) curve:

No matter what $\mu$ and $\sigma$ are, the area between $\mu$-$\sigma$ and $\mu$+$\sigma$ is about 68%; the area between $\mu$-2$\sigma$ and $\mu$+2$\sigma$ is about 95%; and the area between $\mu$-3$\sigma$ and $\mu$+3$\sigma$ is about 99.7%. Almost all values fall within 3 standard deviations.
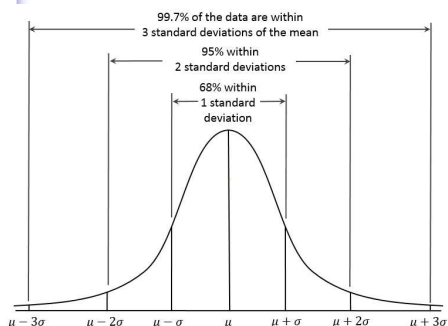
47

## 68-95-99.7 Rule of bell curve



48

## Summary of Symbols

- $S^2$ = Sample variance
- $S$ = Sample standard dev
- $\sigma^2$ = Population (true or theoretical) variance
- $\sigma$ = Population standard dev.
- $\overline{X}$ = Sample mean
- $\mu$ = Population mean
- IQR = interquartile range (middle 50%)

49

## Examples of bad graphics

50



THE SHRINKING FAMILY DOCTOR
In California

Percentage of Doctors Devoted Solely to Family Practice

| 1964 | 1975 | 1990 |
| --- | --- | --- |
| 27% | 16.0% | 12.0% |

1: 4,232
6,212

1: 3,167
6,694

1: 2,247 RATIO TO POPULATION
8,023 Doctors

*Los Angeles Times*, August 5, 1979, p. 3.

**What's wrong with this graph?**

from: ER Tufte. The Visual Display of Quantitative Information. Graphics Press, Cheshire, Connecticut, 1983, p.69

51

From: Visual Revelations: Graphical Tales of Fate and Deception from Napoleon Bonaparte to Ross Perot Wainer, H. 1997, p.29.

**52**



Correctly scaled X-axis…

**53**



**What's wrong with this graph?**

from: ER Tufte. The Visual Display of Quantitative Information. Graphics Press, Cheshire, Connecticut, 1983, p.74

**54**

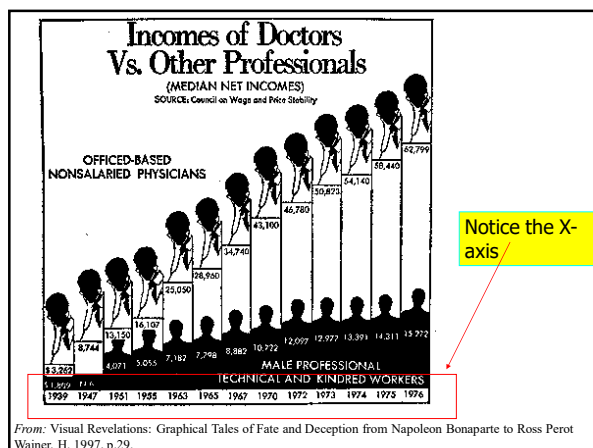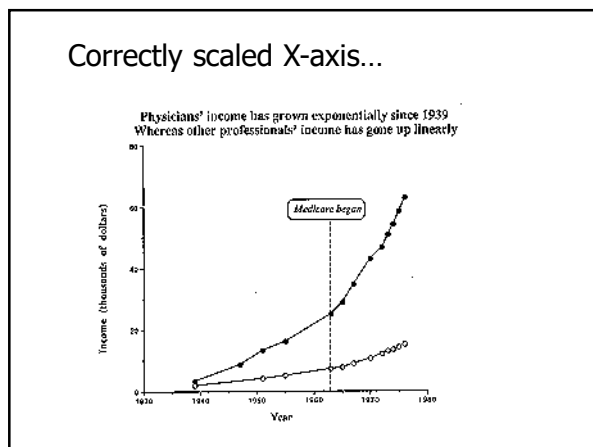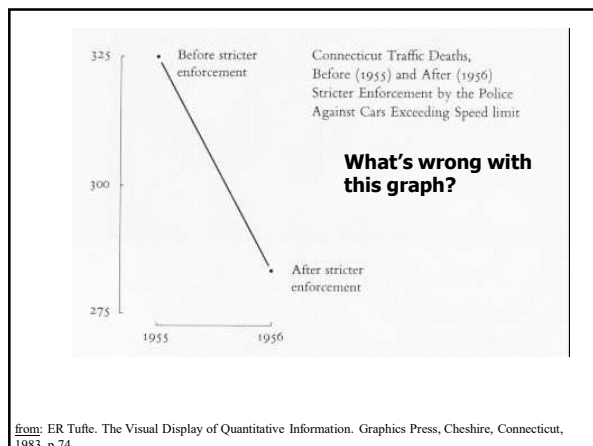**A few more data points add immensely to the account:**

Connecticut Traffic Deaths, 1951–1959

**55**

What's the message here?

**Sotheby's / Christie's**
Worldwide Sales
Market Share Analysis

| 1985 | 1986 | 1987 | 1988 | 1989 | 1990 |

42% 58% | 40% 60% | 41% 59% | 44% 56% | 42% 58% | 44% 56%

SOTHEBY'S
CHRISTIE'S

Market Share Analysis With Buyer's Premium

*Diagraphics II*, 1994

**56**

Sotheby's / Christie's Worldwide Sales

Sotheby's

Christie's

*Diagraphics II*, 1994

**57**

19

**58**



**59**

# References

- http://www.math.yorku.ca/SCS/Gallery/
- Kline et al. *Annals of Emergency Medicine* 2002; 39: 144-152.
- Statistics for Managers Using Microsoft® Excel  4th Edition, 2004 Prentice-Hall
- Tappin, L. (1994). "Analyzing data relating to the Challenger disaster". *Mathematics Teacher*, 87, 423-426
- Tufte. The Visual Display of Quantitative Information. Graphics Press, Cheshire, Connecticut, 1983.
- Visual Revelations: Graphical Tales of Fate and Deception from Napoleon Bonaparte to Ross Perot Wainer, H. 1997.

**60**