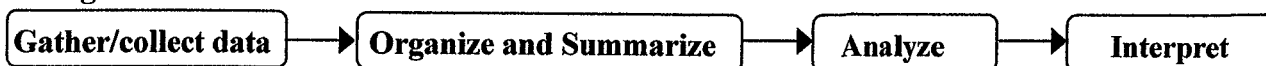


Chapter 1: Introduction to Statistics

STATISTICS IS THE STUDY OF DATA

Things we do with data:



Definitions:

Population

Sample

Variable: the description in words of the characteristic of interest

Data: the information collected about the variable for individuals in the population or sample.

- *The variable is a sentence explaining the question you are asking in order to obtain information,*
- *The data are the information obtained as answers to the question.*

After the data is collected, we use the data to calculate number that summarizes and describe the characteristic for the population or sample. This number is called a parameter or a statistic.

Parameter

Statistic

Example 1: Suppose we are interested in the ages of De Anza College (DAC) students

Population:

Sample:

Variable:

Data:

Examples of Parameters:

- **POPULATION PROPORTION:** In a recent quarter, 39% of all DAC students were over age 25
- **POPULATION AVERAGE:** In a recent quarter, the average age of all DAC students was 27.1 years
- **POPULATION MEDIAN:** In a recent quarter, half of all DAC students were age 22 years or less

Examples of Statistics

- **SAMPLE PROPORTION:** In a recent quarter, 39% of a sample of DAC students were over age 25
- **SAMPLE AVERAGE:** In a recent quarter, the average age of a sample of DAC students was 27.1 years
- **SAMPLE MEDIAN:** In a recent quarter, half of a sample of DAC students were age 22 years or less

This work is licensed under a [Creative Commons Attribution-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-sa/4.0/).

Introduction to Statistics Notes, by Roberta Bloom De Anza College.

Some material is derived from Introductory Statistics from Open Stax (Illowsky/Dean) available for download for free at <http://cnx.org/content/11562/latest/> or <https://openstax.org/details/introductory-statistics>



CHAPTER 1: TYPES OF VARIABLES AND DATA:

EXAMPLE 2: Suppose we are studying the commutes of De Anza College students to school from home.

What is the population? _____

Give a possible example of a sample:

Some examples of variables and data:

Variable: "the distance that a (one, individual) student commutes to De Anza College "

Data: 2.5 miles, 8.4 miles, 0.25 miles, 52 miles, . .

Variable: "how a De Anza student commutes to school

Data: car, car, bus, bike, walk , car, bus, bus, car , bike, car

Variable: "The number of days per week a student comes to campus

Data: 4, 5, 5, 2, 4, 3, 5, 5, 2, 5, 1, 3, 4, 0, 1, 3, 5, . . .

Quantitative

○ Quantitative Discrete

○ Quantitative Continuous

Qualitative (also called Categorical):

LEVELS OF MEASUREMENT

Nominal: qualitative data: "how a student commutes to school" or "a student's major"

Ordinal: data that can be ranked:

- 1st place, 2nd place, 3rd place, . . . last place, such as in a race or competition
- preference ratings such as
1 = strongly agree, 2 = agree, 3 = indifferent, 4 = disagree, 5 = strongly disagree
where there is no set measurement associated with each of the different responses.

Interval: data has constant differences for intervals

But there is no definite natural "0" value that can be used as the base point for comparisons

Example: Temperature degrees Celsius or degrees Fahrenheit:

These are arbitrary scales; comparisons differ when using °C versus °F

80 degrees F = 26.7 degrees C and 40 degrees F = 4.4 degrees C

Is it twice as hot (80 degrees F is two times 40 degrees F)?

OR is it 6 times as hot (26.7 degrees C is about 6 times 4.4 degrees C)?

Lack of a natural zero base point, ratio comparisons do not make sense.

Ratio: data has constant differences for intervals

AND there is a natural "0" value that can be used as a base point for ratio comparisons

Example: distance travelled from home to school:

Asaf travels 3 miles and Ben travels 6 miles. Ben travels twice as far as Asaf

CHAPTER 1: TYPES OF VARIABLES AND DATA:

EXAMPLE 3: A weight loss clinic is studying the exercise habits of its clients.

A survey of all client records shows that 64% of all clients exercise regularly.

A survey of a random sample of 100 clients who exercise regularly shows that of those clients in the sample, the average amount of time that clients exercise each week is 2.5 hours, and that they exercise on average 4 days each week.

a. For each variable listed below, indicate if it is qualitative, quantitative discrete, or quantitative continuous?

Whether a client exercises regularly: _____

The amount of time a client exercises each week: _____

The number of days a client exercises each week: _____

b The number 64% above is a : A. data B. statistic C. parameter D. sample E. population

c. The numbers 2.5 and 4 above are: A. data B. statistics C. parameters D. sample E. population

d. A client responds that she exercises 3 days a week for 2.4 hours total; the numbers 3 and 2.4 are _____

STATISTICAL METHODS

Descriptive statistics

Inferential statistics

TYPES OF STATISTICAL STUDIES:

A. Census vs Sampling

1. Census:

2. Sampling:

B. Observational Study vs Experiment

1. Observational Study:

2. Experiment:

| | | | |
|---|---------------|----------|--------------|
| Some Vocabulary for Experiments: | Treatment | Placebo | Response |
| | Control Group | Blinding | Double Blind |

EXAMPLE 4: An example of an experiment to test a new drug

- A study is done to test a new drug that is supposed to relieve pain.
- The patients enrolled in the study are randomly divided into two groups.
- One group is given the new drug, called the _____.
- The other group is called the _____ group and is given the _____ instead of the treatment.
- The _____ of both groups is measured to determine if the drug is more effective at relieving pain than not receiving the drug.
- The study is _____ because neither the patient nor the doctor knows who is receiving the new drug and who is receiving the placebo.
- If a patient develops problems, the doctor works with the study administrator who knows who is receiving the drug and who is receiving the placebo.
- The doctor, the study administrator, and a statistician are part of a team of people who evaluate the effectiveness of the drug based on the results of the study.

How this relates to the real world: As patients, we would not want to use a drug or a medical device unless it has been shown to be safe and effective. In the Silicon Valley area, there are many biotech companies developing drugs and medical devices. They need to conduct studies of the safety and effectiveness of these drugs or devices. UC Santa Cruz Extension has a professional "clinical trials" certification program. People enrolling in this program are expected to already have college degrees in related fields and in this program they study topics such as how the law applies to medical studies, the ethics of medical studies, statistics, design of experiments and other topics.

CRITICAL EVALUATION OF STATISTICAL STUDIES AND RESULTS

IT IS VERY IMPORTANT TO THINK CRITICALLY ABOUT THE VALIDITY AND RESULTS OF STATISTICAL STUDIES RATHER THAN BLINDLY BELIEVING THE RESULTS OF ALL STUDIES.

Common Problems in Statistics to beware of:

- Problems with Samples:
 - A sample should be representative of the population.
 - A sample that is not representative of the population is called "biased".
 - Non-response or refusal of subject to participate
 - Self-Selected Samples might be biased
 - Sample Size Issues
- Collecting data or asking questions in a way that influences the response
- Causality: A relationship between two variables does not necessarily imply that one **causes** the other
They may both be affected by some other variable.
- Self-Funded or Self-Interest Studies
- Misleading Use of Data: improperly displayed graphs, incomplete data, lack of context, not enough information given to understand the data
- Confounding: when the effects of multiple factors on a response can not be separated, it becomes difficult or impossible to draw valid conclusions about the effect of each factor.

EXAMPLE 5:

Study I: Employees of a company are randomly divided into two groups. Group A gets classroom training from an instructor who is available to help and answer questions; Group B gets training via online software with an online discussion board available to get help and answers to questions.

Study II: Researchers are studying whether an early retirement age increases the rate of memory problems in senior citizens. A survey of retired senior citizens showed that those who retired earlier tended to have a higher rate of memory problems after retirement than those who retired at an older age.

Study III: 300 randomly selected individuals are asked if they had been on a diet in the last 8 weeks and how much their weight has changed over the last 8 weeks. Weight change for dieters and non-dieters are compared.

Study IV: 100 individuals are put on a low fat diet, 100 on a low carb diet and 100 eat their normal diet. Their weight change over an 8 week period is recorded.

a. For each of the above, what type of study is it: Observational Study or Experiment?

I: _____

II: _____

III: _____

IV: _____

b. What problem can you see in Study II?

c. Which weight loss study (III or IV) do you think would give the best information about the effect of diet on weight loss? In class, discuss some reasons why.

EXAMPLE 6: A large city is proposing a parcel tax to support education. Each property owner would be assessed a tax of \$100 per property per year. The parcel tax will be voted on by voters in the next election. It will pass if 2/3 of the voters vote in favor of the tax.

I. A group of parents and teachers supporting the parcel tax randomly select and call residents in the city. They identify themselves as members of the Parent Teachers Association for the school system and ask the person who answers the telephone call if they support the parcel tax.

II. A TV news station in the city conducts a "Facebook" survey. Viewers are asked whether they favor or oppose the tax and are given instructions to visit the TV stations Facebook page to respond about their opinion. The poll is publicized and responses are solicited by announcements on the TV station's evening news programs.

III. A professional polling organization conducts a survey by randomly calling selected residents in the city. If the resident is a registered voter, he or she is asked his/her their opinion about the proposed parcel tax. They are asked whether they favor the tax, oppose the tax, or have no opinion. These three choices are presented to the individual in random order, so that not all respondents hear the choices in the same order.

a. Which survey do you think would produce the most accurate prediction of the election results?

b. For each of the other two surveys, what problems do you think there might be with the information obtained? Discuss some of the problems with these surveys in class.

CHAPTER 1: TYPES OF SAMPLES

A sample is a part of or a subset of a population.

A sample should be representative of the population

A sample that is not representative of the population is called "biased".

Vocabulary and Concepts:

Sampling Error: Random error obtained by using part of the population to represent the whole population

Non-Sampling Error: Non-random error: improper data collection recording or sampling techniques, bias,

Random Sample: Every person has equal chance of being included in the sample

Sampling Methods:

- **simple random sample**

- **systematic sample**

- **stratified sample**

- **cluster sample**

- **convenience sample**

EXAMPLE 7: (*a-f from Example 1.6 Collaborative Statistics, Chapter 1, B. Illowsky & S. Dean; free at www.cnx.org*)
Determine the type of sampling method used:

- To form a recreational soccer team, a soccer coach randomly selects 6 players from a group of boys ages 8 to 10, 7 players from a group of boys ages 11 to 12, and 3 players from a group of boys age 13 to 14. _____
- For a survey of human resource (HR) personnel at high tech companies, a pollster interviews all HR personnel in each of 5 randomly selected high tech companies. _____
- A medical researcher for a hospital interviews every third cancer patient from a list of cancer patients at that local hospital. _____
- A high school counselor uses a computer to generate 50 random numbers and then selects students whose names correspond to the numbers. _____
- A student interviews classmates in his algebra class to determine how many pairs of jeans a student owns, on average. _____
- In a study to learn what types of after school child care are used in their district, a school district administrator randomly selects 6 classes at each school and surveys all parents with children in the selected classes. _____

CHAPTER 1: ORGANIZING DATA

NUMERICAL SUMMARIES OF QUALITATIVE DATA: Tables and Graphs – see worksheet

NUMERICAL SUMMARIES OF QUANTITATIVE DATA: FREQUENCY DISTRIBUTIONS

Quantitative (numerical) data can be summarized in a frequency distribution table.

The data may be presented individually (ungrouped) or grouped into intervals.

- A frequency distribution counts the number of data items that fall into each interval.
Frequency = count = number of data values that lie in the interval.
- A relative frequency distribution shows the proportion (fraction or percent) of data items in each interval.
Relative Frequency = proportion of data values that lie in the interval = $\frac{\text{Frequency}}{\text{Number of Observations}}$.
- **Cumulative Relative Frequency = sum of relative frequencies up to and including current interval or value**

EXAMPLE 8:
Individual Data Values
(ungrouped data)

Plants are being studied in a lab experiment. For a sample of 16 plants, the number of flowers on a plant are:
 2,5,3,1,2,4,1,2,3,1,1,2,7,4,2,3

| X = Number of Flowers | Frequency | Relative Frequency | Cumulative Relative Frequency |
|-----------------------|-----------|--------------------|-------------------------------|
| 1 | | | |
| 2 | | | |
| 3 | | | |
| 4 | | | |
| 5 | | | |
| 6 | | | |
| 7 | | | |

- a. What percent of plants had exactly 3 flowers?
- b. What percent of plants had at most 3 flowers?
- c. What percent of plants had more than 3 flowers?
- d. What percent of plants had at least 5 flowers?

Data Values Grouped into Intervals

Lengths of songs, in seconds, for a sample of 65 songs played on a streaming music channel.

What percent of songs were at most 4 minutes long?

What percent of songs were more than 4.5 minutes long?

| X = Song Length, in seconds | Frequency | Relative Frequency | Cumulative Relative Frequency |
|-----------------------------|-----------|--------------------|-------------------------------|
| 166-180 | 3 | 0.046 | 0.046 |
| 181-195 | 6 | 0.092 | 0.138 |
| 196-210 | 7 | 0.108 | 0.246 |
| 211-225 | 12 | 0.185 | 0.431 |
| 226-240 | 7 | 0.108 | 0.539 |
| 241-255 | 14 | 0.215 | 0.754 |
| 256-270 | 4 | 0.062 | 0.816 |
| 271-285 | 6 | 0.092 | 0.908 |
| 286-300 | 3 | 0.046 | 0.954 |
| 301-315 | 1 | 0.015 | 0.969 |
| 316-330 | 2 | 0.031 | 1 |