

Module 3: Data Collection

3.1. Experimental and Observational Studies

Observational study

In an **observational study**, the researcher observes values of the response variable for the sampled subjects, but nothing is done to the subjects themselves (i.e., no treatment is imposed)

Experimental study

In an **experimental study**, or an experiment, a researcher assigns subjects to experimental conditions (or **treatments**). The researcher *manipulates* the treatment variable in order to determine if this *causes* a change in the response.

Only properly designed experimental studies can identify cause and effect relationships.

Observational studies vs. experiments

It is not possible to establish cause-and-effect relationships with observational studies because of *extraneous variables*.

An **extraneous variable** is any variable that potentially influences the variable being studied (known as the *response variable*).

A **confounding** variable is a variable associated with the treatment or explanatory variable.

Well-designed experiments statistically eliminate the effect of confounding variables. Because treatments are randomly assigned, groups of subjects receiving different treatments will be balanced with respect to extraneous variables. For example, groups will be similar in terms of age, lifestyle, gender, and so on.

Does smoking cause cancer?

How do we know that smoking causes cancer? Although lung cancer rates are higher for smokers than non-smokers, we must consider and be able to dismiss several important confounding variables.

- Drinking: Drinking and smoking habits are correlated. How do we know that it is smoking, rather than drinking, that causes lung cancer?
- Exercise: Non-smokers are more likely than smokers to exercise regularly. Perhaps exercise can prevent cancer and is negatively correlated with smoking.
- Eating habits: It is likely that smokers do not eat as many vegetables as non-smokers. If eating vegetables prevents cancer, this might explain the association seen between smokers and lung cancer.
- Genetics: Is there an aspect of somebody's genetic makeup that makes them more likely to both smoke and to develop lung cancer?

Now consider the following experiment to determine whether smoking (tobacco) might cause cancer:

One hundred genetically identical mice are divided into two groups of 50 mice each. The first group is exposed to a fixed amount of tobacco smoke 3 times a day for a year; the second group was exposed only to air in the room. All mice were put on the same diet.

There are still extraneous variables that influence cancer rates. However, because this is an *experiment* these variables will not differ (statistically) between the treatment groups.

1. Genetics and diet are **controlled**.
2. All other factors (including exercise) are **randomized**.
3. The **treatment** is exposure to tobacco smoke

Any difference in cancer rates between the two groups must be due to tobacco smoke exposure, or random chance.

3.2: Observational studies and sampling

When observational studies are appropriate

Not all studies seek to identify causal relationships between two variables.

For example, observational studies are appropriate for

- estimating the average height of all adult females in the United States
- gauging public opinion on a topic, such as capital punishment
- marketing studies to assess how people rate a new product

The goal of an observational study (and an experiment) is to draw conclusions about a larger population of interest. This means that

1. the sample must be *representative* of the population
2. the variables being studied must reflect the question of interest; in a survey, participants must answer honestly and accurately

Example

Suppose I want to know what percent of Eastern students prefer Google to Bing. One possible sampling design would be to sample the students in this class. Is this class representative of the population of all Eastern students?

No! Our sample should be representative of all Eastern students with respect to gender, age, class status, demographic data, proportion of student-athletes, etc. This is an example of a sampling design chosen out of convenience but one that does not represent the population

The key to selecting a sample that is representative of the population is to sample in such a way that each individual has the same chance of being in the sample. This is known as a **simple random sample**.

A **convenience sample** is a sample chosen out of convenience and is almost never representative of the population.

How survey questions are worded is extremely important

Consider these survey questions (from 2003):

- Do you favor or oppose taking military action in Iraq to end Saddam Husseins rule?
- Do you favor or oppose taking military action in Iraq to end Saddam Husseins rule, even if it meant that U.S. forces might suffer thousands of casualties?

Despite the fact that the above questions are 'essentially' asking the *same thing*, 68% favor military action in the first case while only 43% favor military action in the second case

Source: <http://www.pewresearch.org/methodology/u-s-survey-research/questionnaire-design/>

3.3: Summary

Summary

- One can only conclude that there is a cause-and-effect relationship based on data collected from a well-designed experiment.
- The sample used to conduct the study must be representative of the population of interest.
- For survey questions, the wording must be carefully chosen; when measurements are taken, these must be well-defined and carefully considered

All of these should be taken into account when considering limitations or weaknesses of a research project.

3.4: Examples

Student Drug Testing Not Effective in Reducing Drug Use

A drug use questionnaire was filled out by over 76,000 students in 497 high schools and 225 middle schools nationwide. This study led to the headline 'Student Drug Testing Not Effective in Reducing Drug Use'.

Drug Tests?	Drug Use		n
	Yes	No	
Yes	0.37	0.63	5,653
No	0.36	0.64	17,437

1. Is this an observational study or an experiment? **An observational study since no variables are being manipulated. An experiment would have randomly selected the schools that perform drug testing.**

2. Do you agree with the headline 'Student Drug Testing Not Effective in Reducing Drug Use'? No! It is possible that drug testing was performed primarily in schools with high drug use rates to begin with. Even if this led to a decrease in drug use, then overall you might not see a difference in drug use between schools that do and do not administer drug tests.

Which diet is better?

A group of 100 individuals are divided into 2 groups as follows:

Diet	Diet #1	Diet #2
Doctor	Dr. Jones	Dr. Smith
Exercise	jogging for 20 minutes, 3x weekly	jogging for 20 minutes, 3x weekly
Time	3 months	3 months
Weight loss	18 lbs	11 lbs

The researchers conclude that Diet #1 is better. Is this conclusion correct?

Is this the correct conclusion?

No! The doctor is completely confounded with the the treatment! It is possible that the personality of the doctor, rather than the diet itself, could explain the difference between the two groups. One solution is to use the same doctor for both groups, ideally in a way that the doctor does not know which patient is receiving which treatment.