## Section 2.5 - The Question of Causation

Statistics 104

Autumn 2004



Copyright ©2004 by Mark E. Irwin

## Causation

- Does smoking cause cancer?
- Did chemical weapons exposure cause health problems in Gulf War vets?
- Will increasing the speed limit increase traffic fatalities?
- Will bringing storks into an area increase the birth rate?

Example:

Brothers and sisters heights are highly correlated. However a tall brother doesn't cause a tall sister.

A more likely cause: common genetics

Even though there may no be a causal relationship between two variables, it can still be useful to predict a sister's height from her brother's height.

A lack of association doesn't mean that there is no causal relationship. An important lurking variable could have been missing from the analysis.

Strong relationships could mean a number of things.

1. Causation (x causes y)

High temperatures in the summer lead to higher electricity use (fans, air conditioning, etc)



(Note: dashed line indicates association, arrow indicated causation)

2. Common response

An unobserved variable leads to changes in x and y.

It has been observed that children with more cavities tend to have larger vocabularies.

However it is hard to see how more cavities might lead to larger vocabularies (or vice versa).



However in this case, both variables are associated with age.

3. Confounding

The effect x has on y is mixed up with other variables

Example: Strength of molded parts

x: time in mold

y: strength of part

In study, higher strength was associated with longer mold times.

The way the experiment was performed was to have all the samples at 10 seconds in the mold done first, then the samples at 20 seconds, then 30 seconds, and so on.

They also saw a strong relationship between strength and the order done. The time in the mold and the order done were confounded.

It ended up that the mold got warmer as more batches were done and higher temperature increases strength.



Examples:

1. Education and salaries in Economics

A study was performed which showed a mean salary ordering of

PhD < MA < BA

However in private industry

 $\mathsf{PhD} > \mathsf{MA} > \mathsf{BA}$ 

In Government, the same pattern

PhD > MA > BA

Teaching salaries tend to be less than in industry and government.

Where a person chooses to work is confounded with education level.

How can this happen?

Company A		Company B		
n	$ar{x}$	n	$\overline{x}$	
50	\$16,400	50	\$16,000	

Now lets break things down by length of employment

	Company A		Company B	
Employed	n	$ar{x}$	n	$ar{x}$
< 5 years	10	\$10,000	40	\$12,500
> 5 years	40	\$18,000	10	\$20,000

In this hypothetical example, Company A tends to have many more long term employees, which leads to a higher overall average salary, even though for both experience levels, Company B pays better. 2. Car Weight and Gas Mileage

In this data set I believe there are all three situations: Causation, Common Response, and Confounding.

- Causation: Physics says the more weight, the more energy you need to move an object, therefore implying worse gas mileage.
- Common response: the type of car (van, sports, SUV, etc) influences the weight, plus other factors that affect the gas mileage of the car.
- Confounding: While weight has a causative effect, its actually effect can not be accurately ascertained since weight is counfounded with a number of factors, such as engine size or horsepower.

## **Establishing Causation**

Best is by experiment

**Experiment:** A designed study where levels of explanatory variables are set by the experimenter and lurking variables can be controlled.

Lets think about designing an experiment to examine whether smoking causes cancer.

- Choose people at birth to be smokers or non-smokers. Can control for race, gender, etc.
- Set the level of smoking 0,  $\frac{1}{2}$ , 1, or 2 packs per day.
- Observe people until they are 50 and see how the rates differ among the different smoking levels.

Of course you can't do this, one for practicality reason, but more importantly, for ethical reasons.

While the above was presented in a humourous fashion (at least that was the intent), this is related to a serious problem of today.

An active area of research is how to show whether a treatment against a biochemical agent, such as anthrax, ebola, sarin, etc, is effective. It is not possible to expose anybody to these agents.

The current approach is to give the treatment to people, without exposure to the harmful agent to investigate the safety of the treatment. To examine the effectiveness of the treatment, an experiment is run on animals where some are given the treatment, some aren't, and the response after exposure to the agent is compared.

So where possible, experiments are best for establishing causation, but other approaches such as observational studies have to be used in other cases.

In the case where experimentation is possible, repeated experimental evidence is best.

## FDA requirements for a new drug

**Stage 1** Doseage and toxicology studies

Stage 2 Small scale screening and efficacy trials

**Stage 3** At least 2 large scale randomized controlled trials

Stage 4 Followup and monitoring studies

If only observational studies are available

- The association is strong
- The association is consistent
- The association must hold when plausible other variables (possible confounders) are taken into account

- Dose response: increasing dose is associated with stronger response for example
- Time ordering: the alleged cause precedes the effect in time
- Plausible explanation of causal link must exist.

Types of observational studies:

- Retrospective: subjects are questioned about past events
  - Smoking habits
  - Occupation
  - Past health problems
- Prospective: a study group is chosen and followed over time
  - Yearly examinations
  - Checks for changes, in smoking, occupation, etc

A famous example of a prospective study is the Framingham Heart Study, which began in 1948.

- **Original Cohort** 5209 men and women aged 30 to 62 from Framingham MA were entered into the study in 1948. As of February 1998, 1095 were known to be alive.
- **Offspring Cohort** 5135 men and women who are offspring of the original cohort and their spouses was established. As of February 1998, 4524 were still alive, with 20 lost to followup and 4 in whom survival status was unknown.
- **Generation III Cohort** In 2002, a third cohort was being recruited. The goal is to enroll approximately 3500 grandchildren of the original cohort to understand how genetic factors relate to cardiovascular disease.

Where possible, prospective studies are preferred to retrospective studies since the quality of data is usually superior.

However retrospective studies can suggest future prospective studies or experiments.