Economics Research Methods

Peter Howley

Professor of Behavioral Economics

Economics Division

Leeds University Business School
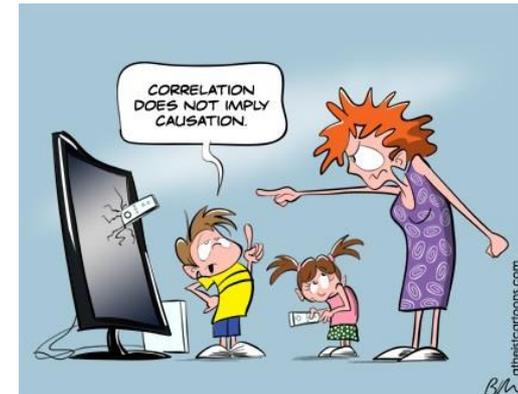
Web: https://business.leeds.ac.uk/divisions-economics/staff/126/dr-peter-howley

# Correlation v causality: Thinking critically

- You have all heard "correlation does not cause causation"
    - but we get this wrong a lot!

- How can factors be correlated but not causally related?

- Pure chance but also

    - **omitted variable bias** (third variable problem)

    - **bi-directional causality** (reverse causality)

*Economists often refer to this as endogeneity bias*

- At the end of the lecture you should be able to explain in your own words (using examples) what these terms mean
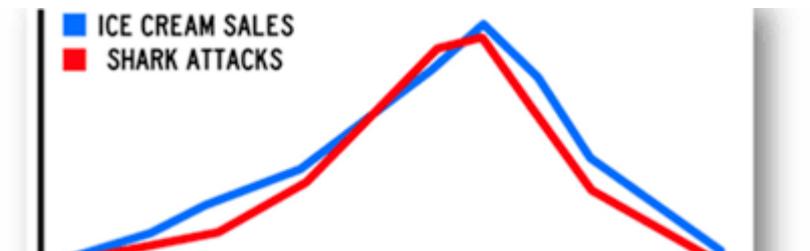
# Correlation v causality: Thinking critically

- Correlation: Degree of relationship between two variables
  (e.g. don't attend versus attend lectures) and grade-point average

- Correlation coefficient indicates the strength of the relationship between two variables (e.g., r = 0.50)

**Where the problems occur**

- Many many things are significantly related with each other – but a **<u>significant relationship</u>** does not mean **<u>a causal one</u>**

# The magic p-value!

- **Hypothesis testing**: - statement that you want to test. In general, the Null hypothesis ($H_0$) is that things are the same as each other, or the same as a theoretical expectation.
  - e.g. no **relationship** between attendance and exam performance class, males earn the same as females, no relationship between smoking and cancer

- When you perform a hypothesis test in statistics, a **p-value** helps you determine the significance of your results, i.e. whether you can reject your Null Hypothesis …

- The **p-value** is a number between 0 and 1 and interpreted in the following way: A small **p-value** (typically $\leq 0.05$) indicates strong evidence against the null hypothesis, so you can reject it and accept the alternative hypothesis ($H_1$).

  - put differently the lower the p-value the more likely it is that the relationship you observe is statistically significant, or in other words
  - A low p-value means that it is unlikely you would obtain the observed results if the Null hypothesis were true

# P-values and hypothesis testing

- **Type 1 error** also commonly referred to as a false positive: Rejection of a true Null Hypothesis

- **Type 2** error also commonly referred to as a false negative: Non-rejection of a false null hypothesis

- Both are very common (Type 1 generally considered a more serious problem)

    - can never eliminate risk but can take constructive steps to minimise occurrence (we will discuss this throughout the module)

# The magic p value

- Null Hypothesis: ($H_0$) No relationship between attendance at class and exam performance

  - Alternative hypothesis: ($H_1$) Attendance at class does lead to better exam performance

- How can I test this – collect (analyse) all data relating to attendance and exam results. Correlate class attendance with exam performance (observational study)

- Result: Correlation coefficient of 0.40 and p-value of 0.04
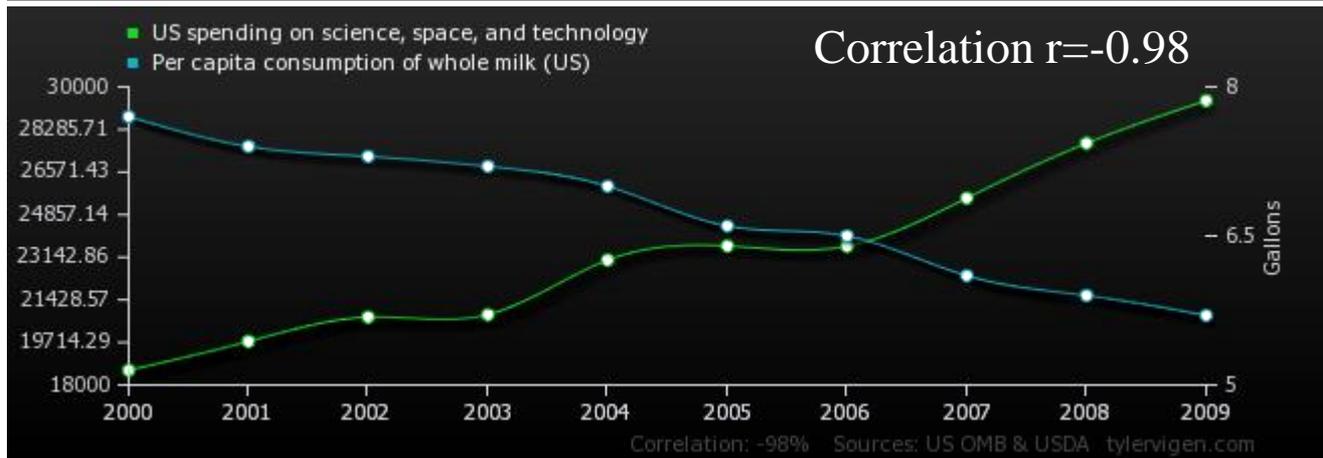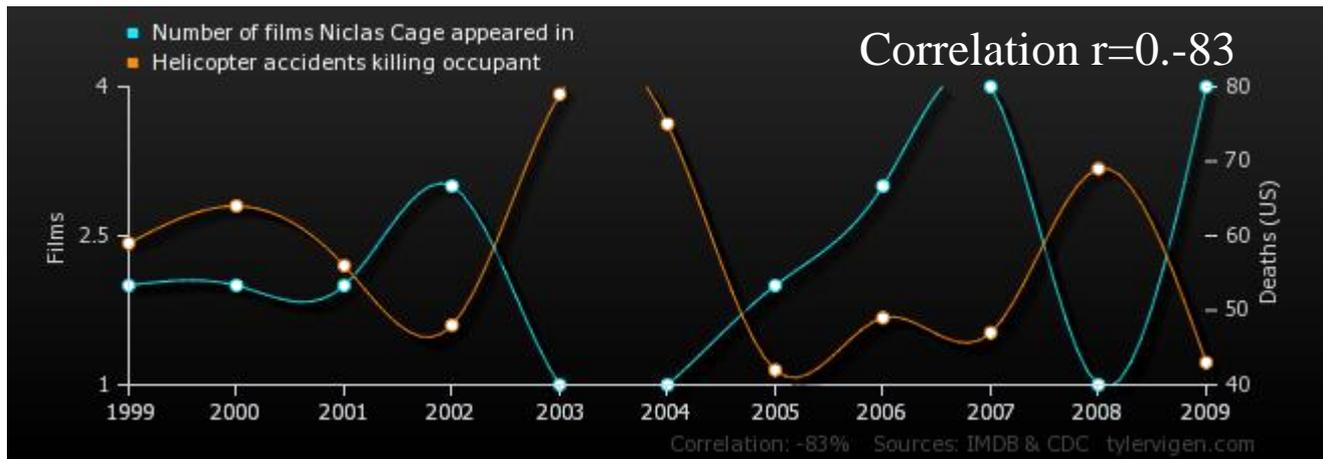
- Interpret this for me?

# The magic p value

- Null Hypothesis: ($H_0$) Attendance at class leads to better examine performance
  - Alternative hypothesis: ($H_1$) Attendance at class does not lead to better exam performance

- How can I test this – collect all data relating to attendance and exam results. Correlate class attendance with exam performance

- Result: Correlation coefficient of 0.40 and p-value of 0.04

- Interpret this for me?

- Can we conclude that attendance at class has a positive impact on exam performance?

# The magic p-value

- **A conclusion does not immediately become 'true' on one side of the divide and 'false' on the other."**

- **P-values are a helpful guide when it comes to examining the relationship between variables, but you need to do much more than rely on p-values when it comes to causality**

- **Don't underestimate the importance of logic/reason and or theory when developing your hypothesis and interpreting your results**

    **- just as important as p values!**

- **How could variables be significantly related with each other but not causally related?**

# 1: Pure chance



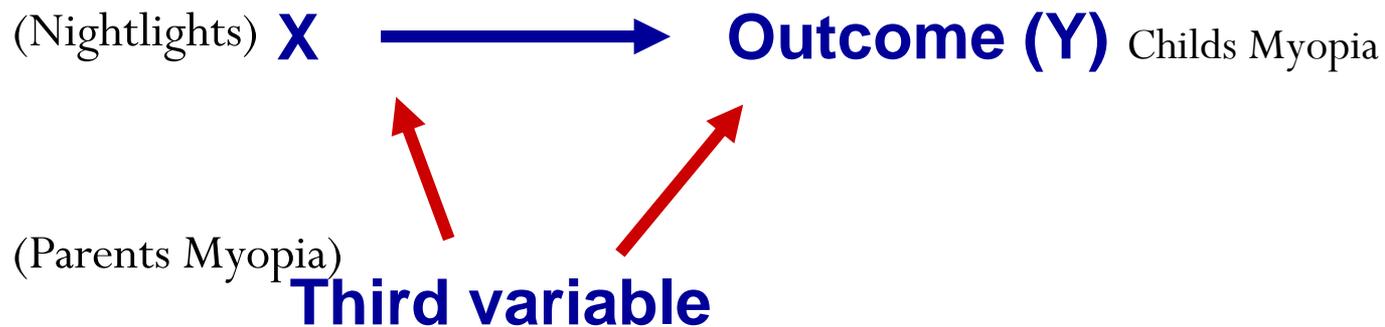Correlation r=0.-83

Correlation r=-0.98

http://t.co/vWOyN0N1IB

# 2. Omitted variable bias (also commonly referred to as confounding – example from Monday)
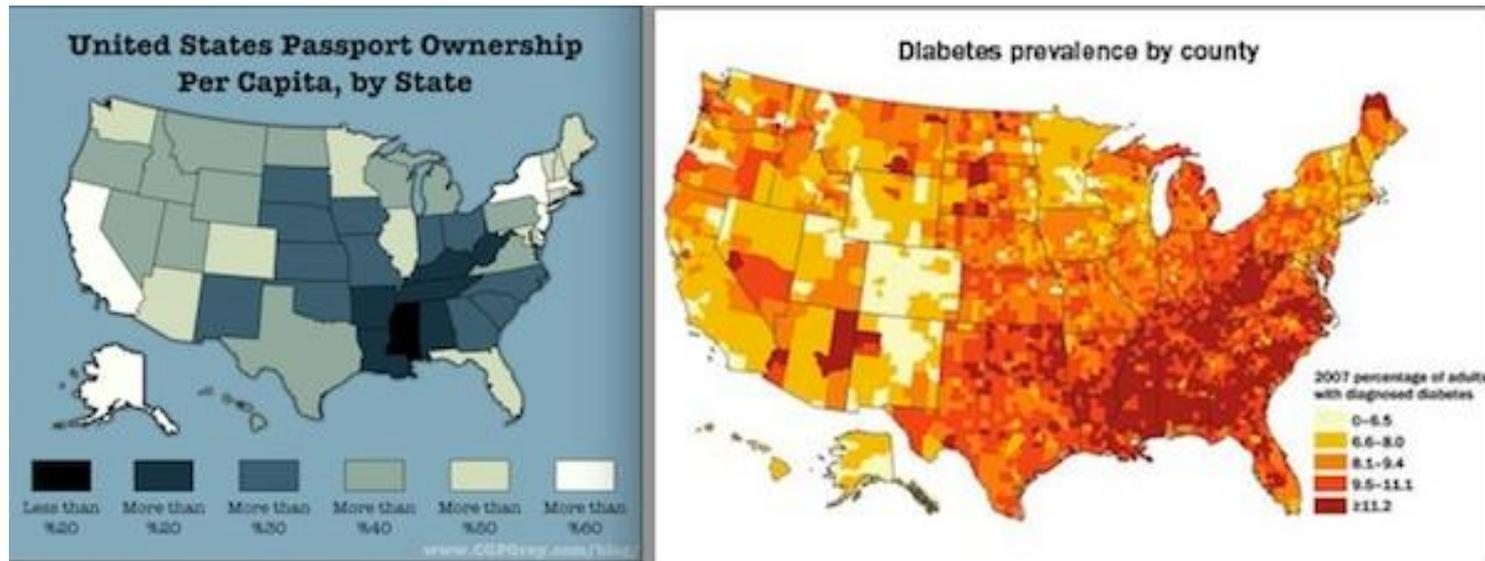
- Third variable may account for the association (**omitted variable bias**)

- What does this mean?

 X is correlated with Y, but this correlation is because both X and Y are correlated with a third variable (Z)
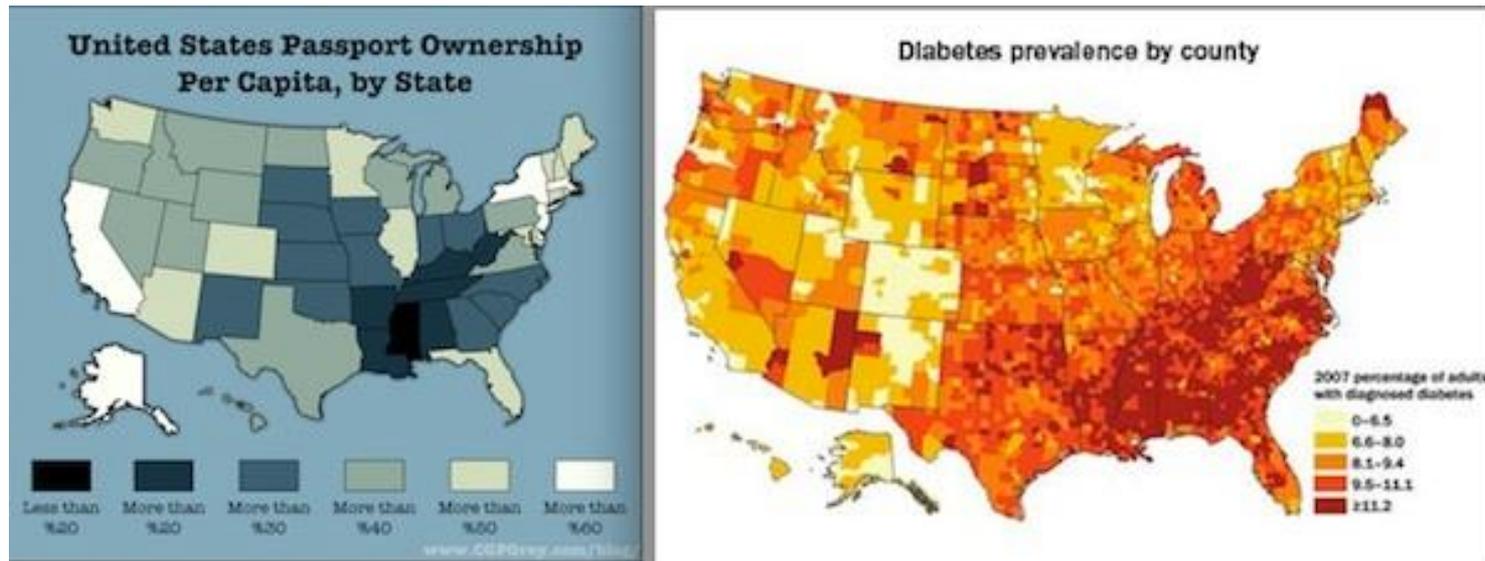
(Nightlights) **X** $\longrightarrow$ **Outcome (Y)** Childs Myopia

(Parents Myopia) **Third variable**

# Lets look at some further examples

- People who slept with their shoes on are very likely to wake up with a headache.
    - ?Therefore, sleeping with shoes on causes headache

# Lets look at some examples

- People who slept with their shoes on are very likely to wake up with a headache.

    ?Therefore, sleeping with shoes on causes headache



Okay above examples are somewhat obvious – well at least I hope so!
We will look at less obvious examples that you are more likely to see published or reported by the media but have the same problem over the coming weeks

# More serious examples

- More seriously, when hormone replacement therapy became commonplace, doctors noticed that women taking HRT seemed less likely to get coronary heart disease. Some doctors suggested a causal relationship - that HRT lowered the risk of heart disease.

# More serious examples

- More seriously, when hormone replacement therapy became commonplace, doctors noticed that women taking HRT seemed less likely to get coronary heart disease. Some doctors suggested a causal relationship - that HRT lowered the risk of heart disease.

- Again it turned out that there was a third variable at play (*omitted variable bias*). Women who were taking HRT were more likely to come from higher socio-economic groups, with healthier diet and exercise habits. It's this that lowered the risk of heart disease. In the end, other tests showed that HRT actually raised the risk slightly.

- These omitted variables meant there was a type of **selection bias**

# Consider this – truckloads of research findings based on this premise frequently reported in the popular press

**Imagine this study**: *Parents were asked to indicate which of the following foods their children ate at least twice a month*:

| | | |
|---|---|---|
| Apple Pie | Fried Chicken | Okra |
| Baked Potatoes | Grape Jelly | Peanut Butter |
| Beets | Hamburger | Pizza |
| Broccoli | Hot Dogs | Popcorn |
| Carrots | Hummus | Potato Chips |
| Chicken Soup | Ice Cream | Strawberry Jam |
| Chocolate Cake | Mac and Cheese | Sushi |
| Corn | Mashed Potatoes | Tacos |
| Eclairs | M & Ms | Tomatoes |
| French Fries | Nachos | Twinkies |

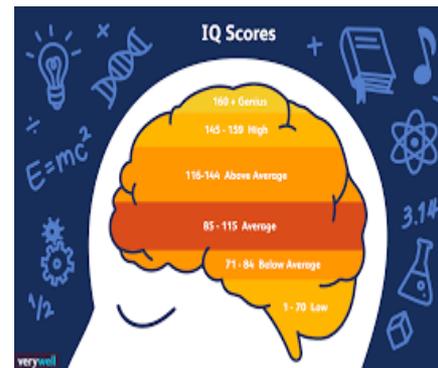Test scores of the children were then examined.

# Eating Sushi Makes You Smarter?

*Imagine this headline*: Scientists reported this week that children who eat sushi score higher on vocabulary tests than children who don't (p < 0.05). The effect of other foods was also studied, but statistically significant results were obtained only for sushi. For example, peanut butter did not show this effect.

Dieticians at local schools, after being informed of the results, said that they will add sushi to their school lunch program.

*Ok, so study suggests a strong relationship between eating sushi and test scores: Can you make a causal inference here?*
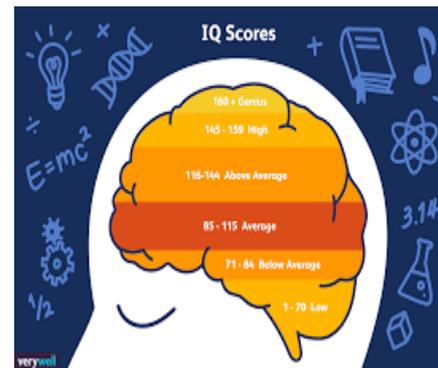
 =  ?

# Eating Sushi Makes You Smarter?

*Imagine this headline*: Scientists reported this week that children who eat sushi score higher on vocabulary tests than children who don't (p < 0.05). The effect of other foods was also studied, but statistically significant results were obtained only for sushi. For example, peanut butter did not show this effect.

Dieticians at local schools, after being informed of the results, said that they will add sushi to their school lunch program.

*Ok, so study suggests a strong relationship between eating sushi and test scores: Can you make a causal inference here? Lots of omitted variable bias (e.g. family income)*

# Selection bias (self-selection)

- *Selection bias leading to a problem of omitted variable bias*: When individuals select themselves into a study (e.g. **choose to eat sushi**) they may not be representative of the population leading to a biased comparison (may differ on unobservable characteristics – omitted variable bias)

- One further common example relates to arguments about the impact of public v private school education on future earnings/test scores

  - A clear gap exists, **but**

  - is the gap the result of better/different education or is it simply that people who choose (or chosen for them) to go to fee-paying schools differ on a host of other attributes that are likely related to exam performance?

18

# 3. Reverse (bi-directional) causality

What does this mean

- X is correlated with Y, but perhaps Y is correlated with X!

# Lets look at some examples

- Previous studies have linked drinking diet soda with a higher risk of diabetes, among other health problems.

- ? Therefore, no more

- Remember example from last week: Walking quickly can prevent depression

# Always have **a critical eye** when Interpreting statistical results

Examples of headlines from newspaper articles – article implied a causal relationship

- *People who sleep 6 hours a night live longer then eight hours or more.*

- *Relationship between Parental Restrictions on Movies and Adolescent Use of Tobacco and Alcohol*

- *Passive smoking dents children's IQ?*

- *Preschoolers more likely to become bullies if they watch lots of TV*

All these were based on a simple correlation between two groups of people.

**Can easily construct a narrative to imply a causal relationship – and this is the problem**

# Policy implications

- Real policy discussions are often infected by spurious implications of causality, often offered disingenuously


- When you see scientists and the media reporting a connection between x and y,' it's really important to be critical about whether there's a causal mechanism."

- One of the golden rules of statistics is that **correlation does not equal causation**. Just because the movements of two variables track each other closely over time (and you obtain the magic p-value!) doesn't mean that one causes the other.


- Some very humorous examples:
http://www.tylervigen.com/spurious-correlations

# Revision for next lecture

Things I want you to be able to do based on this lecture

Explain in your words (with examples) the following terms

- **Omitted variable bias (confounding – third variable problem)**
- **How does selection bias arise – use an example to make sure you understand the intuition**

- **Bi-directional causality**

- **What a p-value tells you and doesn't tell you!**
- Have an intuitive understanding of the various pitfalls when relying on correlations/associations when examining the direct relationship between two variables

# Interesting blogs on this topic

- Read/listen to before class next week:

- Correlation v causation: https://theconversation.com/clearing-up-confusion-between-correlation-and-causation-30761

- Selection bias: https://theconversation.com/coronavirus-country-comparisons-are-pointless-unless-we-account-for-these-biases-in-testing-135464

Correlation, causation and confusion:
https://www.thenewatlantis.com/publications/correlation-causation-and-confusion

Misleading statistics in the news: https://www.youtube.com/watch?v=mJ63-bQc9Xg