# Repeat this several times...(we will in this course)
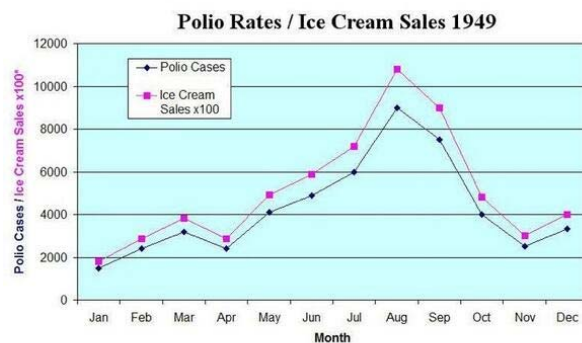
# Funny, not so funny



- Is this crazy? Yes
- But, how about this one…

## Let's repeat it again…



Link: https://www.youtube.com/watch?v=8B271L3NtAw

---

# Part III:
# Selection problem and random assignment

Daniel Schwartz P.

## Part III

- The selection problem
- Random assignment and causal inference
- Treatment effects
- Experiments
- Attrition and sample size

- In other words…
  - How does random assignment solve the selection problem? How do we draw causal inference? How can we run experiments? How can we distinguish good and bad experiments? (learn our strengths and weaknesses…training) How can we analyze results based on experimental work?

## Main references for this Part

- Gerber, A. S., & Green, D. P. (2012). Field experiments: Design, analysis, and interpretation. WW Norton.
- Angrist, J. & Pischke, J-S. (2009) "Mostly Harmless Econometrics: An empiricist companion". Princeton University Press.
- Angrist, J. D., & Pischke, J. S. (2014). Mastering Metrics: The Path from Cause to Effect. Princeton University Press.
- Hoyle, R. H., Harris, M. J., & Judd, C. M. (2002). Research methods in social relations. Thomson Learning.
- Some definitions were based on the course "Experimental Design for Behavioral and Social Sciences" by Howard Seltman
- Wheelan, C. (2013). Naked statistics: stripping the dread from the data. WW Norton & Company. (Chapter 12)
- Bram, U. (2014). Thinking Statistically. Kuri Books.
- Ellis, P. D. (2010). The essential guide to effect sizes: Statistical power, meta-analysis, and the interpretation of research results. Cambridge University Press.

## Motivation

- "The most credible and influential research designs use random assignment" ("Mostly Harmless Econometrics", Angrist & Pischke, p.11)

- "we will eventually see that, when carefully applied, econometric methods can simulate a ceteris paribus experiment." ("Introductory Econometrics, Wooldridge, p.13)

- "In cases where a controlled experiment is impractical or immoral, we need to find some other ways of approximating the counterfactual" ("Naked Statistics", Wheelan, p.240)

… but wait a second, but many times we can run a controlled experiment!

(e.g. class size on education, vaccines on health, marketing campaigns on sales, incentives on blood donation, information on reducing energy consumptions, etc!)

(Why do econometrics books, in general, do not focus on this?)

---

## Application: "Do not kill people with your research" (from Wheelan, 2013)



By DAVID W FREEMAN / CBS NEWS / April 6, 2011, 9:43 AM

**Estrogen pills: Heart disease, cancer risks overblown?**

Source: http://www.cbsnews.com/news/estrogen-pills-heart-disease-cancer-risks-overblown/

## Application 2: "Did eBay Just Prove That Paid Search Ads Don't Work?" (Harvard Business Review)



Paper: Blake, T., Nosko, C., & Tadelis, S. (2015). Consumer Heterogeneity and Paid Search Effectiveness: A Large-Scale Field Experiment. *Econometrica*, *83*(1), 155-174.

Source: https://www.theguardian.com/technology/2013/mar/13/google-keyword-advertising-wastes-money-ebay

---

## The selection problem: Intuition

- Samples may be very useful
  - What happens when they are not truly random?
- The real problems arise when samples are biased and we don't account for that



From "Thinking Statistically" (Uri Bram)

- Think about surveys or feedback systems, which in general are based on voluntary feedback
- How about teacher evaluation?

## The selection problem

- Let's see an example (from Angrist & Pischke, 2011):

  - Let's imagine we are studying a poor elderly population that uses hospitals emergency rooms for primary care. Do this elderly get healthier?

  - This can actually be very problematic

  - One approach: compare health status between people who have been to the hospital and those who haven't.

    - Data from National Health Interview Survey (NHIS) 2005: "during the past 12 months, was the respondent a patient in a hospital overnight? [Health status: from 1 (poor health) to 5 (excellent health)]

| Group | Sample Size | Mean Health Status | Std. Error |
|---|---|---|---|
| Hospital | 7,774 | 3.21 | 0.014 |
| No hospital | 90,049 | 3.93 | 0.003 |

---

## The selection problem

- Let's describe this problem:

  - Think about hospital treatment as a binary random variable: $D_i = \{0,1\}$ and the Outcome of interest (here health status) is denoted by $Y_i$

  - Is $Y_i$ affected by hospital care?

  - In an ideal world we can have this situation:

## Average treatment effect

- We are interested in the average treatment effect (ATE)

- We want: $Y_i(1) - Y_i(0) = \tau_i$
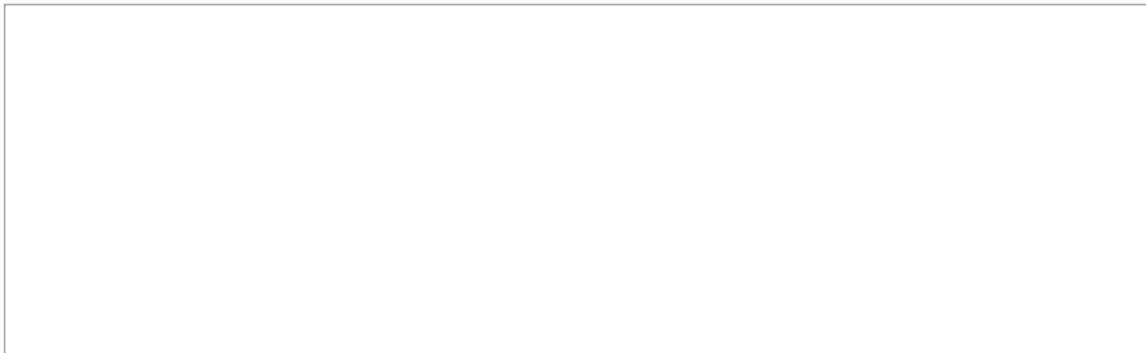  - where $\tau_i$ is the causal effect of the treatment

$$ATE = \frac{1}{N} \sum_{i=1}^{N} \tau_i$$

- But, we can't observe both $Y_i(1)$ and $Y_i(0)$
- When treatments are randomly assigned:

$$ATE = E[Y_i(1)|D_i=1] - E[Y_i(0)|D_i=0]$$

## The selection problem

- But, we don't observe both potential outcomes
  - In the review session you will see an illustration with potential outcomes

- We can compare the average (causal) effect comparing those who were and those who were not hospitalized:

## How can random assignment solve the selection problem?

- Random assignment makes $D_i$ independent of potential outcomes ($Y_i$)
  - Then:
  
    $E[Y_{0i}|D_i=1] = E[Y_{0i}|D_i=0]$, or an equivalent notation:

    $$E[Y_i(0)|D_i=1] = E[Y_i(0)|D_i=0] = E[Y_i(0)]$$

- And therefore, the selection bias is eliminated
- This is great – although we will see how problems may arise in some experiments

## Application with health insurance

- We would like to know whether having health insurance improve health (application from "Mastering' Metrics," Angrist & Pischke)
  - Why is this important?
  - We cannot see someone, at the same time, being insured and not being insured, i.e. we cannot observe both $Y_i(1)$ and $Y_i(0)$
- DV: Self-report health status, from poor to excellent
- Treatment: being insured (control: not being insured)
- Can we compare people who has insurance with those without it? Why not?
  - Actually the former group reports much higher health status
  - E.g. more educated people tend to be healthier

# Application with health insurance (cont'd)

- Let's layout the causal effect using this "imaginary" table
  - Causal effect? What's the selection bias?



| Outcomes and treatments for Khuzdar and Maria | | |
| --- | --- | --- |
| | Khuzdar Khalat | Maria Moreño |
| Potential outcome without insurance: $Y_{0i}$ | 3 | 5 |
| Potential outcome with insurance: $Y_{1i}$ | 4 | 5 |
| Treatment (insurance status chosen): $D_i$ | 1 | 0 |
| Actual health outcome: $Y_i$ | 4 | 5 |
| Treatment effect: $Y_{1i} - Y_{0i}$ | 1 | 0 |

  - What if education is the only difference between groups?
  - What if we toss a coin whether K or M get insurance? (remember the law of the large numbers)
- How can we solve the selection bias?
  - E.g. "The Oregon Trail"