

Doing Applied Research

MIXTAPE TRACK



Section 2. Practical tips for writing your applied paper
Style, tone, length and grammar



- What type of journal are you targeting?
 - Economics and (for instance) medical journals have distinctive styles



Media attention and Vaccine Hesitancy: Examining the mediating effects of Fear of COVID-19 and the moderating role of Trust in leadership

Lulin Zhou, Sabina Ampon-Wireko, Xinglong Xu, Prince Edwudzie Quansah, and Ebenezer Larnyo

Published: February 18, 2022 <https://doi.org/10.1371/journal.pone.0263610>



The current study, therefore, hypothesized that; *H1: Media attention has a significant influence on COVID-19 vaccine hesitancy. H2: Media attention significantly influences fear of COVID-19 positively.*



I see manuscripts submitted to the *Journal of Policy Analysis and Management* (JPAM) with this type of phrasing all the time. I imagine that some policy journals publish papers that use this type of phrasing, but JPAM, as a general rule, does not.

Written in third person



The current study, therefore, hypothesized that; *H1: Media attention has a significant influence on COVID-19 vaccine hesitancy. H2: Media attention significantly influences fear of COVID-19 positively.*



I see manuscripts submitted to the *Journal of Policy Analysis and Management* (JPAM) with this type of phrasing all the time. I imagine that some policy journals publish papers that use this type of phrasing, but JPAM, as a general rule, does not.

The current study, therefore, hypothesized that; H1: *Media attention has a significant influence on COVID-19 vaccine hesitancy.* H2: *Media attention significantly influences fear of COVID-19 positively.*

Conflates statistical significance and importance



Style, Tone, Length and Grammar

- What type of journal are you targeting?
 - Economics and (for instance) medical journals have distinctive styles
- Let's imagine that you're targeting an economics journal
 - Within economics, journals have distinctive styles
 - Read a few articles in the journal that you're targeting
 - *Economics Letters, American Economic Review: Insights, AEA: Papers and Proceedings*, have strict word count limits



Style, Tone, Length and Grammar

- What type of journal are you targeting?
 - Economics and (for instance) medical journals have distinctive styles
- Let's imagine that you're targeting a field journal such as the *Journal of Health Economics*, *Labour Economics*, or the *Journal of Public Economics*



Style and Tone

- Engaging without being chatty or informal
 - More formal than the language you hear being used in seminars
 - Avoid “get” and “got”
 - Avoid using contractions (e.g., won’t, don’t, can’t, shouldn’t)
 - Avoid this construction:

“...we use for our main calibration the value 0.38, i.e. the median of the elasticity estimates in the literature.”
- Read what you wrote out loud
 - Does each sentence set up the next?
 - Does each paragraph set up the next?
- Don’t reinvent the wheel
 - Read some classics. How did Charles and Guryan (2008) begin their introduction? What are they doing in their first paragraph? In their second paragraph?



Prejudice and Wages: An Empirical Assessment of Becker's *The Economics of Discrimination*

Kerwin Kofi Charles and Jonathan Guryan

University of Chicago and National Bureau of Economic Research



I. Introduction

Becker's (1957) seminal *The Economics of Discrimination* launched the formal analysis of labor market discrimination among economists. Becker's analysis focused on the relationship between racial prejudice among whites and discrimination against racial minorities in a competitive model. In contrast to much of the contemporaneous literature, Becker formalized the definition of racial preferences...



I. Introduction

Becker's (1957) seminal *The Economics of Discrimination* launched the formal analysis of labor market discrimination among economists. Becker's analysis focused on the relationship between racial prejudice among whites and discrimination against racial minorities in a competitive model. In contrast to much of the contemporaneous literature, Becker formalized the definition of racial preferences...

In the short-run version of Becker's employer discrimination model, racial prejudice causes some employers to regard black workers as more expensive than they truly are. Market pressures cause blacks to be hired by the least prejudiced employers in the market and to sort away from those with the highest levels of prejudice...Racial wage gaps...will generally be determined by variation in the level of prejudice of those in the lower tail of the prejudice distribution; how prejudiced the most prejudiced employers are should not matter at all for wages in Becker's framework. Finally,...equilibrium wages for blacks should vary negatively with the number of blacks in the market...



Prejudice and Wages: An Empirical Assessment of Becker's *The Economics of Discrimination*

Kerwin Kofi Charles and Jonathan Guryan

University of Chicago and National Bureau of Economic Research

- I. Introduction
- II. Background
- III. Data Summary
- IV. Empirical Strategy
- V. Results
- VI. Conclusion

I'd call this the “classic” way to structure an economics paper



Length

- Aim for 18-23 pages of text, 5-6 figures and tables (combined total)
- More figures than tables
- Extra tables and figures can go in an appendix
- The revision process is usually going to add text, tables, and figures
- Thirty pages of text and 10 figures/tables is too much! Referees and editors have short attention spans.



Grammar

- Pay someone to read you paper over for typos and grammatical errors
- Make sure there is a one-to-one correspondence between in-text citations and references
- Use a consistent style for the references.
 - Include the first names of the authors. The editor is looking at your references for potential referees.

Doing Applied Research

MIXTAPE TRACK



Section 2. Practical tips for writing your applied paper:
Writing the first half of your paper



Writing the first half of your applied paper

The Introduction

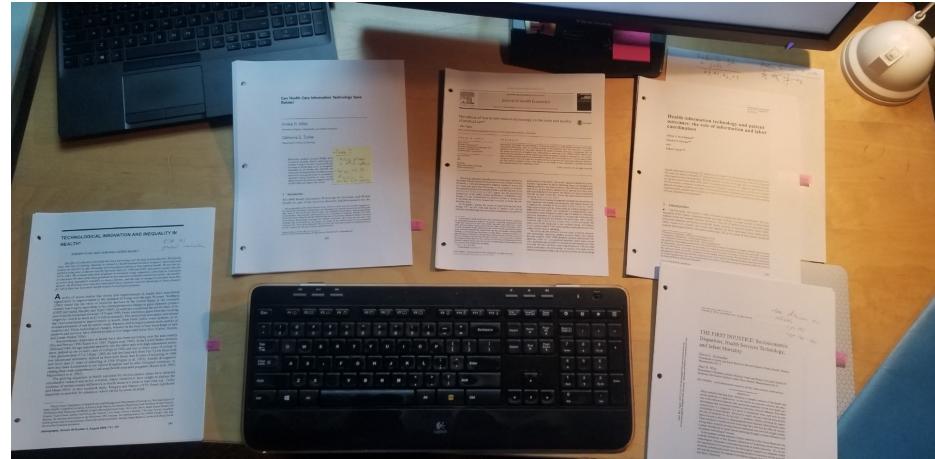
- Perhaps the most important part of your paper.
- Why should someone keep reading your paper if they aren't "hooked" by the end of the introduction?





The Introduction

- When I sat down to write my first introduction, I felt overwhelmed. I did not know where to begin.
- For my first paper, I surrounded myself by well-published papers with well-written introductions...



- I quickly learned that there is a method to writing a good introduction.
 - Just like there is a template for a paper in its entirety, there is a template for the introduction.
- Let's break this thing down...



The Introduction

- An introduction to an applied research paper should contain the following elements:
 - (i) A general motivation. Ideally, a motivation that captivates economists and non-economists alike.
 - (ii) A more specific motivation. For instance, a motivation that appeals to the very policy being analyzed or specific theory being tested.
 - (iii) A brief description of the most important background information or institutional details (save the rest for the background section of your paper)
 - (iv) A description of your data.
 - (v) A description of your identification strategy.
 - (vi) A description of your principal results.
 - (vii) A roadmap for the rest of your paper (optional)
- I like to approach writing my introductions as a process of checking off the above boxes.
- Let's walk through an example...

-You are NOT writing a mystery novel!
-Your reader should know exactly what to expect by the end of the introduction.



Was the First Public Health Campaign Successful?

By D. MARK ANDERSON, KERWIN KOFI CHARLES,
CLAUDIO LAS HERAS OLIVARES, AND DANIEL I. REES*

The US tuberculosis (TB) movement pioneered many of the strategies of modern public health campaigns. Using newly transcribed mortality data at the municipal level for the period 1900–1917, we explore the effectiveness of public health measures championed by the TB movement, including the establishment of sanatoriums and open-air camps, prohibitions on public spitting and common cups, and requirements that local health officials be notified about TB cases. Our results suggest that these and other anti-TB measures can explain, at most, only a small portion of the overall decline in pulmonary TB mortality observed during the period under study. (JEL H51, I12, I18, N31, N32)

In 1900, 194 out of every 100,000 Americans died of tuberculosis (TB), making it the second-leading cause of death, behind only pneumonia/influenza (Jones, Podolsky, and Greene 2012). Although an effective treatment would not be introduced until after World War II (Daniel 2006), the TB mortality rate fell dramatically over the next three decades. By 1920, it had fallen to 113 per 100,000 persons; by 1930, it had fallen to 71 per 100,000 persons (Jones, Podolsky, and Greene 2012).

How was TB vanquished, or at least controlled, in the United States and other developed countries? Scholars have proposed several explanations, including better living conditions, herd immunity due to natural selection, reduced virulence, and improved nutrition (Smith 2003; Daniel 2006; Kunitz 2007, 96–197; Lönnroth et al. 2009; and Mercer 2014, 127–29). The introduction of basic public health measures (e.g., isolating patients in sanatoriums and TB hospitals) is another potential explanation (Wilson 1990, Fairchild and Oppenheimer 1998), but scholars have questioned whether such measures contributed meaningfully to the decline in TB mortality (McKeown 1976, Coker 2003, and Daniel 2006).¹

- (i) A general motivation that appeals to a broad audience.
- (ii) A more specific motivation that appeals to scholars working in this area.
 - Previous explanations for the TB mortality decline are listed
 - Highlight the explanation we are interested in testing (i.e., the introduction of public health measures)
 - Highlight that this remains an open question, which also highlights our scientific contribution

questioned whether such measures contributed meaningfully to the decline in TB mortality (McKeown 1976, Coker 2003, and Daniel 2006).¹

Drawing on newly transcribed data from a variety of primary sources, the current study explores whether the TB movement contributed to the decline in TB mortality in the United States. The movement began with the establishment of the Pennsylvania Society for the Prevention of Tuberculosis in 1892 and gained momentum when the National Association for the Study and Prevention of Tuberculosis (NASPT) was founded in 1904 (Shryock 1957, 52; Teller 1988, 30). Spearheaded by voluntary associations and supported by the sale of Christmas Seals, the US TB movement pioneered many of the strategies of modern public health campaigns (Teller 1988, 1, 121–26; Jones and Greene 2013; and Rosen 1993, 226–31).

Between 1900 and 1917, hundreds of state and local TB associations sprung up across the United States (NASPT 1916, Knopf 1922). These associations distributed educational materials and provided financial support to sanatoriums and TB hospitals, where patients with active TB were isolated from the general population and, if lucky, could recover. In addition, these associations advocated, often successfully, for the passage of legislation designed to curb the transmission of TB, including bans on public spitting and requirements that doctors notify local public health officials about active TB cases.

Although remarkable in its scope and intensity, the effectiveness of the US TB movement has, to date, not been studied in a systematic fashion.² Using municipal-level data for the period 1900–1917 from *Mortality Statistics*, which was published on an annual basis by the US Census Bureau, we estimate the relationship between pulmonary TB mortality and the introduction of public health measures designed to curb the spread of the disease. Our estimates, which control flexibly for common shocks and municipal-level heterogeneity, suggest that most anti-TB measures had no discernable impact on pulmonary TB mortality. Two exceptions stand out: there is evidence, albeit tentative, that requiring TB cases to be reported to local health officials led to a modest reduction in pulmonary TB mortality; likewise, the opening of a state-run sanatorium is associated with a modest reduction in pulmonary TB mortality. However, these two measures can explain, at most, only a small portion of the overall decline in pulmonary TB mortality during the period 1900–1917.



- (iii) Important background details
 - Historical background
 - Relevance to modern day public health
 - What constituted the anti-TB campaign?
- (iv) Data description
- (v) Identification
- (vi) Results



The Introduction: Two Common FAQs

- How many contributions should I list?

- Don't get carried away here. A good paper probably only makes one or two *real* contributions.
- If you feel that your paper really does make more than this, clearly highlight what you think is the most important of the listed contributions.
- Big mistake is to list a bunch of contributions that are not ranked in order of importance and claim to advance a number of different fields.
 - This can send the signal that you really aren't sure where your major contributions lie.
 - Is my paper a labor paper or a behavioral paper or a econ history paper...or wait, is it something else? Try to avoid creating this confusion (especially for editors who are trying to figure out the appropriate referees for your manuscript).

- How long should my introduction be?

- Clearly, there is no rule of thumb.
 - But, avoid lumping your intro and background sections into one big introduction.
 - Erring on the side of brevity is probably better. Do not ask the reader to wade through a 10-page introduction.
- The last 15 papers I have written as full-length articles for economics journals have had introductions of the following lengths (size 12 font, double-spaced):

1-2 pages	2-3 pages	3-4 pages	4-5 pages	5+ pages
3 papers	6 papers	3 papers	3 papers	0 papers!!!



Background

- Break into subsections (not necessarily in this order):
 - (i) Background/historical/institutional details regarding your right-hand side variable of interest
 - Include information in this section that may be important for your identification strategy
 - For instance, how many policy changes do you observe? For how many units (e.g., cities, counties, states, etc.) do you observe pre- and post-treatment data?
 - (ii) Background on your specific relationship of interest
 - (iii) Background specific to the theory you are testing (if any)
 - (iv) Literature review can be a separate section or woven into the above sections.
 - As I have become more comfortable with writing background sections, I like to cite the relevant literature within sections (i)-(iii). Generally flows better.
 - As a matter of style, one should include subsection headers for each of these background components. This helps to quickly orient a reader who may be going back and forth between sections of your manuscript.



Example 1. “The Effects of Pretrial Detention on Conviction, Future Crime, and Employment: Evidence from Randomly Assigned Judges”
(Dobbie et al. 2018, *American Economic Review*)

- **Goal of paper:** Use detention tendencies of quasi-randomly assigned bail judges to estimate the causal effects of pretrial detention on subsequent defendant outcomes.
- **Background, Section A.** Overview of the bail system in the United States
 - Very general overview providing reader with how the bail system works in the United States
- **Background, Section B.** Description of the setting specific to their data (i.e., Philadelphia and Miami-Dade Counties)
 - Broken into 3 subsections: (i) One on the bail system in Philadelphia County; (ii) one on the bail system in Miami-Dade County; (iii) and one describing how the features of the Philadelphia and Miami-Dade bail systems make them an appropriate setting for their research design.
 - I really like their subsection B. (iii). Great example of how one can describe the intuition behind their identification strategy before you get to any regression equations. I *highly recommend* reading this part of their paper.



Example 2. “Personal Belief Exemptions for School-Entry Vaccinations, Vaccination Rates, and Academic Achievement” (Hair et al. 2018, *Journal of Health Economics*)

- **Goal of paper:** Examine how state-level nonmedical exemptions from school-entry vaccine mandates affect vaccination rates in early childhood and academic achievement in middle school.
- Nice example of a straightforward (and well-written) paper that does not require a lengthy background section.
- **Background, Section 1.** Description of policy changes in their treated states (i.e., Texas and Arkansas).
- **Background, Section 2.** Discussion of how exemptions might affect a downstream outcome such as academic achievement.
 - The authors list and describe the possible mechanisms.
 - Great example of a discussion on mechanisms that comes before results are presented.



Example 3. “Safeguarding Consumers Through Minimum Quality Standards: Milk Inspections and Urban Mortality, 1880-1910” (Anderson et al. 2022, NBER Working Paper No. 30063)

- **Goal of paper:** Estimate the effect of municipal-level milk inspections on infant and waterborne mortality.
- **Background, Section 1.** The general history of municipal milk inspections in the United States.
 - Describe exactly how inspections worked in practice.
- **Background, Section 2.** History of other municipal-level efforts to provide pure milk.
 - Give readers a sense of what other public health policies and efforts were being mandated.
 - Relevant for the empirical strategy. What are the omitted variables we be concerned about?
- **Background, Section 3.** Literature review on minimum quality standards (MQSs) and their effects.
 - Review the competing empirical literature and describe the theoretical contributions
 - I've listed this paper for this specific subsection. Of the 3 examples, this paper tests a classic economic theory, whereas the others do not. As a result, the literature review must tackle both empirical and theoretical studies.



Background: A common mistake!

- In one's literature review, *not all previous studies should be treated equally.*
- Researchers should provide careful and accurate evaluations of previous studies.
 - Most policymakers have never heard of fixed effects, have no idea what a regression discontinuity is, and do not care whether state-specific time trends were included on the right-hand side.
 - Policymakers and non-economists count on us to communicate which studies should be taken seriously and which should not (or be ignored).
 - For instance, giving studies that rely on cross-sectional policy variation the same weight as those that rely on within-state policy variation comes with the risk of leading policymakers astray.
- For example, researchers interested in the relationship between marijuana use and alcohol consumption all too often declare that the literature is “mixed”, citing Pacula (1998a) and Williams et al. (2004) as evidence of complementarity and Kelly and Resul (2014) and Sabia et al. (2017) as evidence of substitutability.
 - This is a LAZY characterization of the literature.
 - Pacula (1998a) and Williams et al. (2004) relied on cross-sectional policy variation. Their estimates could simply reflect unobserved factors at the state level such as preferences and attitudes.
 - On the other hand, Kelly and Resul (2014) and Sabia et al. (2017) used panel data and exploited well-defined natural experiments...much more credible research design!



Data and empirical strategy

- Two separate sections or just one?
 - Both ways are common within the greater applied economics literature.
 - More and more, I prefer combining these two components into one section.
 - Every now and then you see someone describe their empirical strategy before their data.
 - I'm not a huge fan of this.
 - As a reader, I want to understand the data set before I start thinking about the empirical strategy in detail
- The ordering I have settled on goes something like the following:
 - First, begin by describing your outcome data.
 - Second, describe your primary estimating equation.
 - Lastly, describe the data for your right-hand side variables.
- This can best be illustrated through another example. Let's go back to the milk inspections paper (Anderson et al. 2022, NBER Working Paper No. 30063)...

4. DATA AND EMPIRICAL FRAMEWORK

The city-level mortality data used in our principal analysis come from a wide range of sources. For the pre-1900 period, mortality counts are from annual municipal and state public health records. These records were obtained either through interlibrary loan, the HathiTrust Digital Library, Google Books, or the archives at the National Library of Medicine in Bethesda, Maryland. For instance, mortality data for Cleveland, OH were obtained from the HathiTrust Digital Library (1880-1883), the Cleveland Public Library via interlibrary loan (1884-1885), and the National Library of Medicine Archives (1886-1899).²⁷ For the period 1900-1910, mortality counts come from *Mortality Statistics*, which, beginning in 1900, was published annually by the U.S. Census Bureau.²⁸ Appendix Table A1 lists the cities and years used in our principal analysis and Appendix Table A2 lists the sources from which the pre-1900 mortality data were obtained.²⁹

Figure 2 shows infant deaths (i.e., deaths among children under the age of one) per 100,000 population for the 35 cities in our sample. Population data come from the decennial censuses and are linearly imputed for intercensal years. In 1880, there were 564.2 infant deaths per 100,000 population. By 1910, the infant mortality rate had fallen to 320.5, or 43 percent.

Figure 2 also shows the waterborne mortality rate, defined as deaths due to diarrheal diseases and typhoid per 100,000 population.³⁰ Typhoid deaths are often used by researchers as a proxy for



- First, we describe the data sources for our outcomes of interest (i.e., mortality rates).
 - If you spent considerable time and effort collecting a unique data set, make sure to play that up! Make sure you get credit for impressive data efforts.
 - Describe your data in a manner that also sheds light on where identification is coming from. It should never be a black box ([App. Table A1](#)).
 - List your sources in detail. Make sure the reader knows exactly where and how you obtained your data ([App. Table A2](#)).
- This is also a great place to show descriptive figures, such as trends in your outcomes ([Figure 2](#))
 - This splices things up a bit.
 - Rather than just a boring description of the data, spend some time highlighting interesting trends, differences in means between groups, etc.

²⁷ To take another example, mortality data for Philadelphia, PA were obtained from the HathiTrust Digital Library (1880-1887, 1890, 1892), the National Library of Medicine Archives (1888-1889, 1894-1896, 1899), and Google Books (1891, 1893, 1897-1898).



Purifying municipal water supplies through filtration and chlorination was primarily a post-1900 phenomenon (Anderson et al. 2022a). It is, however, clear from Figure 2 that infant and waterborne mortality rates had begun trending downwards well before the turn of the 20th century. To explore whether the trends shown in Figure 2 can be explained by milk inspections, we estimate the following regression model:

$$(1) \quad MR_{ct} = \alpha_0 + \nu_c + \lambda_t + \sum_{y=-5}^{-2} \pi_y 1(t - T_c^* = y) + \sum_{y=0}^{10} \pi_y 1(t - T_c^* = y) + \mathbf{X}_c \boldsymbol{\beta} + \varepsilon_{ct}$$

where MR_{ct} is the infant (or waterborne) mortality rate in city c and year $t = 1880\ldots1910$.³¹ City fixed effects, ν_c , control for time-invariant determinants of mortality or reporting differences across municipalities; year fixed effects, λ_t , account for common shocks.

The event-year dummies, y , are equal to 1 when the year of observation is $y = -5, \dots, 0, \dots, 10$ years from T_c^* , the year in which city c 's milk supply was first sampled and inspected with the goal of preventing adulteration and skimming. We omit $y = -1$ from equation (1), which normalizes the estimates of π_y to 0 in that year. The $y = -5$ event-year dummy is equal to 1 if t is 5 or more years before T_c^* . Likewise, the $y = 10$ event-year dummy is equal to 1 if t is 10 or more years after T_c^* .

Inspection dates are listed in Appendix Table A1. They are based on newspaper accounts and other sources such as municipal public health reports.³² The vector \mathbf{X}_c includes controls for whether city c filtered its water, treated it with chlorine, or had undertaken a major clean water project.³³ Appendix Table A4 provides descriptive statistics and Appendix C lists the wide variety of sources used to determine the milk inspection and water-related intervention dates.

- Next, describe your estimating equation.
 - If you can, try to use your descriptive figures/statistics as motivation for your regression equation.
 - Don't screw up your subscripts!!! This happens way too often in papers submitted for review.
- Finally, use your regression equation as a jumping off point to discuss the variables on the right-hand side.
 - Start with describing your primary right-hand side variable of interest.
 - After that, discuss the other Xs.
 - For these other controls, I like to list their sources and other pertinent information in footnotes. As we have done here.
 - Often, these Xs are not all that important. Presenting these variables in this manner allows the reader to decide for herself if she wants the source information.
 - If she does, she can read the footnote. If not, she can skip it.
 - It is common for researchers to go on and on describing the controls in X in a data section before the discussion of the empirical strategy.
 - I do not recommend this.
 - Putting this information in the main text can disrupt the flow of the manuscript, especially for the reader who wants to gloss over this information.



Q&A (\approx 10 minutes)
+
Break (\approx 5 minutes)



Appendix Table A1. Municipal Milk Inspection Dates and Data Availability

City and state	Year milk inspections began	Years covered for infant mortality	Years covered for waterborne mortality
Mobile, AL	1905	1889-1910	1889-1894, 1900-1909
San Francisco, CA ^a	1895	1881-1897, 1900-1910	1881-1887, 1889, 1891-1897, 1900-1910
Bridgeport, CT	...	1880-1910	1880-1882, 1887-1896, 1898-1910
Hartford, CT	1883	1880-1910	1880-1882, 1887-1888, 1890-1896, 1898-1909
New Haven, CT	...	1880-1910	1880-1882, 1887-1896, 1898-1910
Waterbury, CT	...	1880-1910	1880-1882, 1887-1896, 1898-1909
Washington, D.C.	1893	1880-1910	1880-1910
Atlanta, GA	1895	1893-1910	1880-1910
Chicago, IL	1893	1880-1910	1880-1897, 1900-1910
Indianapolis, IN	1896	1884-1887, 1895-1896, 1899-1910	1891-1910
New Orleans, LA	1892	1881-1910	1880-1910
Portland, ME	1902	1887-1910	1894-1898, 1900-1909
Baltimore, MD	1894	1880-1910	1880-1890, 1892-1898, 1900-1910
Grand Rapids, MI	1897	1890, 1892-1910	1890-1894, 1896-1910
St. Paul, MN	1887	1885-1910	1885, 1887, 1889-1910
St. Louis, MO	1887	1880-1896, 1900-1910	1880-1896, 1900-1910
Omaha, NE	1904	1891-1910	1892-1893, 1895-1910
Manchester, NH	1884	1883-1910	1883-1909
Camden, NJ	1883	1880-1910	1880-1886, 1900-1909
Elizabeth, NJ	1883	1880-1910	1880-1886, 1900-1909
Hoboken, NJ	1883	1880-1910	1880-1886, 1900-1909
Jersey City, NJ	1883	1880-1910	1880-1886, 1900-1910
Newark, NJ	1883	1880-1910	1880-1886, 1900-1910
Paterson, NJ	1883	1880-1910	1880-1886, 1900-1910
Trenton, NJ	1883	1880-1910	1880-1886, 1900-1909
Buffalo, NY ^b	1883	...	1881-1910
Rochester, NY ^c	1891	1881-1910	...
Cleveland, OH	1888	1880-1910	1880-1910
Dayton, OH	1887	1880, 1890-1891, 1893-1896, 1898-1910	1880-1910
Toledo, OH	1884	1881-1883, 1885, 1898-1910	1882-1886, 1888, 1894-1898, 1900-1910
Philadelphia, PA	1889	1880-1910	1880-1888, 1890-1894, 1896, 1899-1910
Scranton, PA	1891	1887-1897, 1900-1910	1888-1897, 1900-1910
Charleston, SC	1907	1880-1883, 1885-1894, 1897, 1900-1910	1880-1883, 1885-1894, 1896-1897, 1899-1909

[back](#)



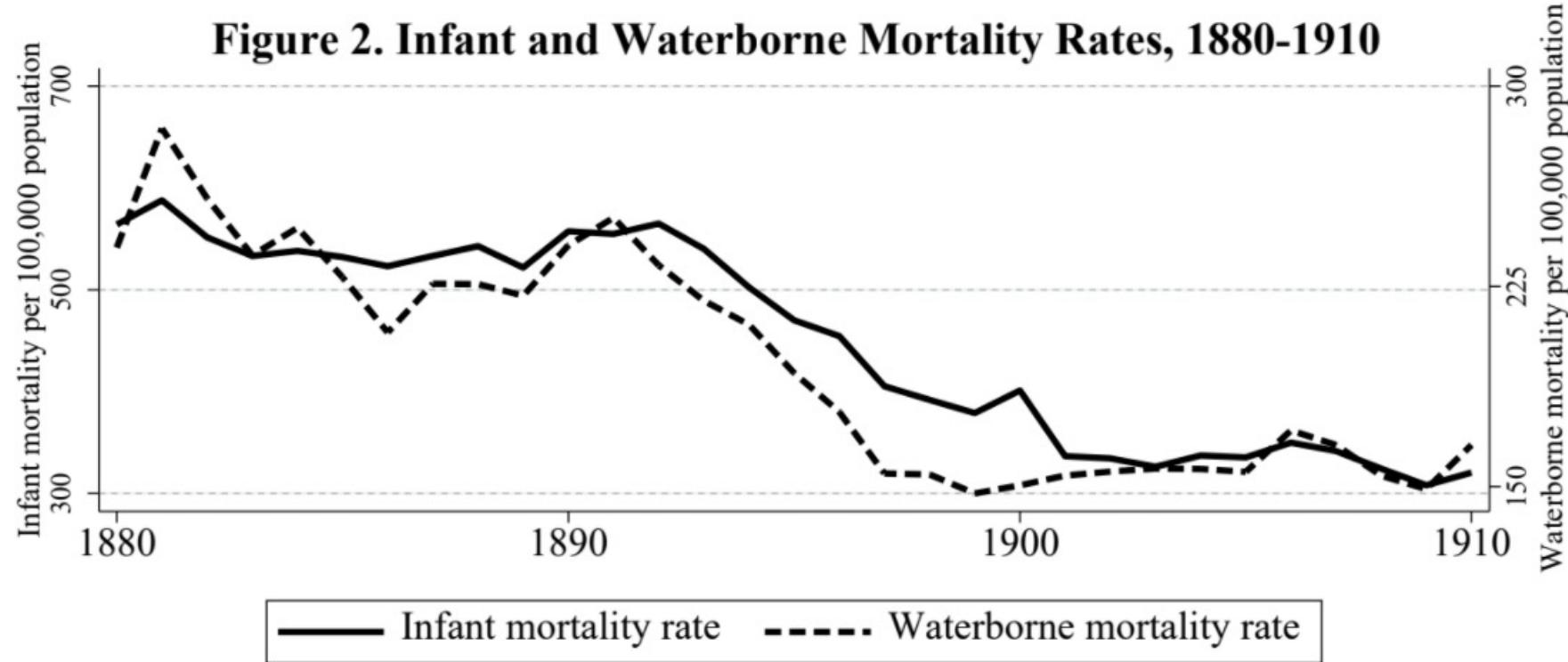
Appendix Table A2. Location of Municipal and State Public Health Reports for Pre-1900 Mortality Data

City and state	Sources
Mobile, AL	National Library of Medicine Archives
San Francisco, CA	National Library of Medicine Archives
Bridgeport, CT	HathiTrust Digital Library
Hartford, CT	HathiTrust Digital Library
New Haven, CT	HathiTrust Digital Library
Waterbury, CT	HathiTrust Digital Library
Washington, D.C.	HathiTrust Digital Library
Atlanta, GA	National Library of Medicine Archives
Chicago, IL	Chicago Municipal Library (1880), HathiTrust Digital Library (1881-1899)
Indianapolis, IN	National Library of Medicine Archives
New Orleans, LA	HathiTrust Digital Library
Portland, ME	National Library of Medicine Archives
Baltimore, MD	HathiTrust Digital Library (1880-1888, 1890, 1892-1899), Google Books (1889), National Library of Medicine Archives (1891), National Library of Medicine Archives
Grand Rapids, MI	National Library of Medicine Archives
St. Paul, MN	National Library of Medicine Archives
St. Louis, MO	St. Louis Public Library (1880-1896), HathiTrust Digital Library (1897-1899)
Omaha, NE	National Library of Medicine Archives
Manchester, NH	HathiTrust Digital Library
Camden, NJ	HathiTrust Digital Library
Elizabeth, NJ	HathiTrust Digital Library
Hoboken, NJ	HathiTrust Digital Library
Jersey City, NJ	HathiTrust Digital Library
Newark, NJ	HathiTrust Digital Library
Paterson, NJ	HathiTrust Digital Library
Trenton, NJ	HathiTrust Digital Library
Buffalo, NY	National Library of Medicine Archives
Rochester, NY	National Library of Medicine Archives
Cleveland, OH	HathiTrust Digital Library (1880-1883), Cleveland Public Library (1884-1885), National Library of Medicine Archives (1886-1899)
Dayton, OH	National Library of Medicine Archives
Toledo, OH	National Library of Medicine Archives
Philadelphia, PA	HathiTrust Digital Library (1880-1887, 1890, 1892), National Library of Medicine Archives (1888-1889, 1894-1896, 1899), Google Books (1891, 1893, 1897-1898)
Scranton, PA	National Library of Medicine Archives
Charleston, SC	National Library of Medicine Archives
Memphis, TN	National Library of Medicine Archives
Richmond, VA	National Library of Medicine Archives
Milwaukee, WI	HathiTrust Digital Library (1880-1884, 1896-1899), University of Wisconsin-Madison

[back](#)



Figure 2. Infant and Waterborne Mortality Rates, 1880-1910



Notes: Based on annual data from municipal and state public health reports (1880-1899) and *Mortality Statistics* (1900-1910).

[back](#)

Doing Applied Research

MIXTAPE TRACK



Section 2. Practical tips for writing your applied paper:

- (i) Producing Effective Tables
- (ii) Useful Phrases When Describing Regression Results
- (iii) Thinking about Statistical Significance and Magnitude
- (iv) Data Mining



Table 1: Estimation results using all cases and variables

	Intercept (μ)		Treatment effect (α)	
	Estimate	Confidence Interval	Estimate	Confidence Interval
own_with_equity	50.88%	(48.63%, 53.13%)	1.12%	(0.33%, 1.92%)
own_no_equity	0.49%	(0.35%, 0.62%)	0.01%	(-0.07%, 0.08%)
rent	43.35%	(41.14%, 45.56%)	-0.87%	(-1.66%, -0.09%)
has_mortgage	60.73%	(58.24%, 63.23%)	1.07%	(0.22%, 1.92%)
home_equity	\$102,585.26	(\$93,756.00, \$111,414.53)	\$2,226.40	(\$324.21, \$4,128.60)
behind	2.49%	(1.91%, 3.07%)	-0.02%	(-0.19%, 0.14%)
months_behind	3.20	(2.24, 4.15)	-0.00	(-0.34, 0.33)
likely_behind	8.67%	(6.94%, 10.40%)	-0.03%	(-0.38%, 0.32%)
foreclose	11.95%	(5.14%, 18.76%)	0.03%	(-2.41%, 2.47%)
housing_asst	3.98%	(3.24%, 4.71%)	0.08%	(-0.12%, 0.28%)
in_public	13.26%	(11.78%, 14.74%)	0.15%	(-0.35%, 0.65%)
gvt_rent	4.93%	(3.77%, 6.09%)	0.14%	(-0.19%, 0.47%)
int_rate	4.59%	(4.43%, 4.74%)	-0.01%	(-0.04%, 0.01%)

Point estimates and 95% confidence intervals of the parameters from (1) as found using imputations generated via the primary OHIE/PSID data fusion model that includes all cases and variables.

This table is from a submission to a top field journal. What can be improved?



The title could be much more informative. Must read the paper to understand what the authors are doing.

Table 1: Estimation results using all cases and variables

	Intercept (μ)		Treatment effect (α)	
	Estimate	Confidence Interval	Estimate	Confidence Interval
own_with_equity	50.88%	(48.63%, 53.13%)	1.12%	(0.33%, 1.92%)
own_no_equity	0.49%	(0.35%, 0.62%)	0.01%	(-0.07%, 0.08%)
rent	43.35%	(41.14%, 45.56%)	-0.87%	(-1.66%, -0.09%)
has_mortgage	60.73%	(58.24%, 63.23%)	1.07%	(0.22%, 1.92%)
home_equity	\$102,585.26	(\$93,756.00, \$111,414.53)	\$2,226.40	(\$324.21, \$4,128.60)
behind	2.49%	(1.91%, 3.07%)	-0.02%	(-0.19%, 0.14%)
months_behind	3.20	(2.24, 4.15)	-0.00	(-0.34, 0.33)
likely_behind	8.67%	(6.94%, 10.40%)	-0.03%	(-0.38%, 0.32%)
foreclose	11.95%	(5.14%, 18.76%)	0.03%	(-2.41%, 2.47%)
housing_asst	3.98%	(3.24%, 4.71%)	0.08%	(-0.12%, 0.28%)
in_public	13.26%	(11.78%, 14.74%)	0.15%	(-0.35%, 0.65%)
gvt_rent	4.93%	(3.77%, 6.09%)	0.14%	(-0.19%, 0.47%)
int_rate	4.59%	(4.43%, 4.74%)	-0.01%	(-0.04%, 0.01%)

Point estimates and 95% confidence intervals of the parameters from (1) as found using imputations generated via the primary OHIE/PSID data fusion model that includes all cases and variables.

Variable names are directly from Stata program

In general, the table is too cluttered. There is too much going on...



Measure in 1000s of \$

Table 1: Estimation results using all cases and variables

	Intercept (μ)		Treatment effect (α)	
	Estimate	Confidence Interval	Estimate	Confidence Interval
own_with_equity	50.88%	(48.63%, 53.13%)	1.12%	(0.33%, 1.92%)
own_no_equity	0.49%	(0.35%, 0.62%)	0.01%	(-0.07%, 0.08%)
rent	43.35%	(41.14%, 45.56%)	-0.87%	(-1.66%, -0.09%)
has_mortgage	60.73%	(58.24%, 63.23%)	1.07%	(0.22%, 1.92%)
home_equity	\$102,585.26	(\$93,756.00, \$111,414.53)	\$2,226.40	(\$324.21, \$4,128.60)
behind	2.49%	(1.91%, 3.07%)	-0.02%	(-0.19%, 0.14%)
months_behind	3.20	(2.24, 4.15)	-0.00	(-0.34, 0.33)
likely_behind	8.67%	(6.94%, 10.40%)	-0.03%	(-0.38%, 0.32%)
foreclose	11.95%	(5.14%, 18.76%)	0.03%	(-2.41%, 2.47%)
housing_asst	3.98%	(3.24%, 4.71%)	0.08%	(-0.12%, 0.28%)
in_public	13.26%	(11.78%, 14.74%)	0.15%	(-0.35%, 0.65%)
gvt_rent	4.93%	(3.77%, 6.09%)	0.14%	(-0.19%, 0.47%)
int_rate	4.59%	(4.43%, 4.74%)	-0.01%	(-0.04%, 0.01%)

Point estimates and 95% confidence intervals of the parameters from (1) as found using imputations generated via the primary OHIE/PSID data fusion model that includes all cases and variables.

What is the treatment? Use the notes to provide some sense of what is being estimated and how.



Period here instead of colon

Table 1. Years of schooling and hourly wages: Sample restricted to Danish respondents

	Men		Women	
	wages	$\ln(\text{wages})$	wages	$\ln(\text{wages})$
<i>Schooling</i>	5.64 (5.11)	.054* (.030)	7.83** (3.55)	.094*** (.007)
Observations	12,385		12,012	

Notes: The sample is restricted to CHAT respondents ages 25-45 from Colorado, Nebraska, and Iowa whose wage was greater than or equal to \$3.25 per hour. Each column reports an estimate from a separate OLS regression of wages on years of schooling. Controls include race, ethnicity, marital status, union status, years of experience, tenure at current job, and firm size. Standard errors clustered at the state level are reported in parentheses.

* Statistically significant at the 10% level; ** at the 5% level; *** at the 1% level.

Now you can use a colon here



Don't report too many digits past the decimal point—it gives the impression that your estimates are precise to the nth digit past the decimal point

Table 1. Estimates of the relationship between years of schooling and hourly wages

	Men		Women	
	wages	ln(wages)	wages	ln(wages)
<i>Schooling</i>	5.64 (5.11)	.054* (.030)	7.83** (3.55)	.094*** (.007)
Observations	12,385		12,012	

Notes: The sample is restricted to CHAT respondents ages 25-45 from Colorado, Nebraska, and Iowa whose wage was greater than or equal to \$3.25 per hour. Each column reports an estimate from a separate OLS regression of wages on years of schooling. Controls include race, ethnicity, marital status, union status, years of experience, tenure at current job, and firm size. Standard errors clustered at the state level are reported in parentheses.

* Statistically significant at the 10% level; ** at the 5% level; *** at the 1% level.

Always tell readers if you're reporting t-statistics or SEs under the estimated coefficients



Keep the table uncluttered. Don't report the estimated coefficients of the controls.

Table 1. Estimates of the relationship between years of schooling and hourly wages

	Males		Females ^a	
	wages	ln(wages)	wages	ln(wages)
<i>Schooling</i>	5.64 (5.11)	.054* (.030)	7.83** (3.55)	.094*** (.007)
Observations	12,385		12,012	
State Fixed Effects	Yes		Yes	
State-Specific Trends	Yes		Yes	

Notes: The sample is restricted to CHAT respondents ages 25-45 from Colorado, Nebraska, and Iowa whose wage was greater than or equal to \$3.25 per hour. Each column reports an estimate from a separate OLS regression of wages on years of schooling. Controls include race, ethnicity, marital status, union status, years of experience, tenure at current job, and firm size. Standard errors clustered at the state level are reported in parentheses.

* Statistically significant at the 10% level; ** at the 5% level; *** at the 1% level.

^a Female respondents under the age of 30 were excluded from the analysis.

You can provide additional information in the notes to the tables or under the estimates.



Table 1. Estimates of the relationship between years of schooling and hourly wages

	Males		Females	
	<i>wages</i>	<i>ln(wages)</i>	<i>wages</i>	<i>ln(wages)</i>
<i>Schooling</i>	5.64	.054*	7.83**	.094***
	(5.11)	(.030)	(3.55)	(.007)
Observations	12,385		12,012	

Notes: The sample is restricted to CHAT respondents ages 25-45 from Colorado, Nebraska, and Iowa whose wage was greater than or equal to \$3.25 per hour. Each column reports an estimate from a separate OLS regression of wages on years of schooling. Controls include race, ethnicity, marital status, union status, years of experience, tenure at current job, and firm size. Standard errors clustered at the state level are reported in parentheses.

* Statistically significant at the 10% level; ** at the 5% level; *** at the 1% level.

Ideally, you provide enough information in the table so that the reader can interpret the estimates without having to read the text.



Useful Phrases When Describing Regression Results

- I find evidence of a negative relationship between years of schooling and earnings.
- The results suggest that schooling is negatively related to earnings.
- Consistent with the hypothesis that educational attainment affects earnings, I find that...
- A one-year increase in schooling is associated with a 6.8 percent increase in earnings.
- The estimate of β is .068, suggesting that an additional year of schooling leads to a 6.8 percent increase in earnings.
- The estimated effect of an additional year of schooling on earnings is positive 6.8 percent.



Common Pitfalls When Describing Regression Results

- Don't focus on statistical significance
- Do your best to characterize magnitude
- Be clear about the units
- Don't confuse “percent” and “percentage”



Table 1. Estimates of the relationship between union status and hourly wages

	Men		Women	
	<i>wages</i>	<i>ln(wages)</i>	<i>wages</i>	<i>ln(wages)</i>
<i>Union</i>	5.64	.156	7.83	.098
	(5.11)	(.030)	(3.55)	(.007)
Observations	12,385		12,012	

Notes: Each column reports an estimate from a separate OLS regression of wages on union status. Controls include race, ethnicity, marital status, union status, years of experience, tenure at current job, and firm size. Standard errors clustered at the state level are reported in parentheses.

Union membership is associated with a \$5.64 increase in the hourly wage.

Your estimate of β is .156. The percent change is $(e^{.156} - 1) \times 100 = 16.88$



Table 1. Estimates of the relationship between years of schooling and the probability of being vaccinated

	<i>Vaccinated</i>	<i>Vaccinated</i>
<i>Years of schooling</i>	.111 (.030)	.021 (.007)
State FEs	no	yes
Mean of the DV	.266	.266
Observations	12,385	12,012

Notes: The sample is restricted to CHAT respondents ages 25-45 from Colorado, Nebraska, and Iowa whose wage was greater than or equal to \$3.25 per hour. Each column reports an estimate from a separate OLS regression of vaccination status (0/1) on years of schooling. Controls include race, ethnicity, marital status, union status, years of experience, tenure at current job, and firm size. Standard errors clustered at the state level are reported in parentheses.

An additional year of schooling is associated with a .021 increase in the probability of being vaccinated (with the state fixed effects)



Table 1. Estimates of the relationship between years of schooling and the probability of being vaccinated

	<i>Vaccinated</i>	<i>Vaccinated</i>
<i>Years of schooling</i>	.111 (.030)	.021 (.007)
State FEs	no	yes
Mean of the DV	.266	.266
Observations	12,385	12,012

Notes: The sample is restricted to CHAT respondents ages 25-45 from Colorado, Nebraska, and Iowa whose wage was greater than or equal to \$3.25 per hour. Each column reports an estimate from a separate OLS regression of vaccination status (0/1) on years of schooling. Controls include race, ethnicity, marital status, union status, years of experience, tenure at current job, and firm size. Standard errors clustered at the state level are reported in parentheses.

An additional year of schooling is associated with a 2.1 percentage point increase in the likelihood of being vaccinated (with the state fixed effects)



Table 1. Estimates of the relationship between years of schooling and the probability of being vaccinated

	<i>Vaccinated</i>	<i>Vaccinated</i>
<i>Years of schooling</i>	.111 (.030)	.021 (.007)
State FEs	no	yes
Mean of the DV	.266	.266
Observations	12,385	12,012

Notes: The sample is restricted to CHAT respondents ages 25-45 from Colorado, Nebraska, and Iowa whose wage was greater than or equal to \$3.25 per hour. Each column reports an estimate from a separate OLS regression of vaccination status (0/1) on years of schooling. Controls include race, ethnicity, marital status, union status, years of experience, tenure at current job, and firm size. Standard errors clustered at the state level are reported in parentheses.

Try to include mean of the DV when you can

An additional year of schooling is associated with a 2.1 percentage point increase in the likelihood of being vaccinated, or an 8 percent increase relative to the sample mean ($.021/.266 = .079$)



Table 1. Estimates of the relationship between the unemployment rate and the probability of being vaccinated

Don't report "0.000"		
	<i>Vaccinated</i>	<i>Vaccinated</i>
<i>Unemployment rate</i>	.00005 (.0006)	.005 (.002)
State FEs	no	yes
Mean of the DV	.266	.266
Observations	12,385	12,012

Notes: Each column reports an estimate from a separate OLS regression of vaccination status (0/1) on years of schooling . Controls include race, ethnicity, marital status, union status, years of experience, tenure at current job, and firm size. Standard errors clustered at the state level are reported in parentheses.

An a one percentage point increase in the unemployment rate is associated with a .005 increase in the probability of being vaccinated (with the state fixed effects).

If the unemployment rate goes from 10.1 to 11.1 percent, that is a one percentage point increase in the unemployment rate.



Stargazing

Estimates of the relationship between years of schooling and hourly wages

	Males		Females	
	<i>wages</i>	<i>ln(wages)</i>	<i>wages</i>	<i>ln(wages)</i>
<i>Schooling</i>	5.64	.054 *	7.83 **	.094 ***
	(5.11)	(.030)	(3.55)	(.007)
Observations	12,385		12,012	

Notes: The sample is restricted to CHAT respondents ages 25-45 from Colorado, Nebraska, and Iowa whose wage was greater than or equal to \$3.25 per hour. Each column reports an estimate from a separate OLS regression of wages on years of schooling. Controls include race, ethnicity, marital status, union status, years of experience, tenure at current job, and firm size. Standard errors clustered at the state level are reported in parentheses.

* Statistically significant at the 10% level; ** at the 5% level; *** at the 1% level.



[Journal of Economic Literature](#) > [Vol. 34, No. 1, Mar., 1996](#) > The Standard Error o...



The Standard Error of Regressions

Deirdre N. McCloskey and Stephen T. Ziliak

Journal of Economic Literature

Vol. 34, No. 1 (Mar., 1996), pp. 97-114

Published by: [American Economic Association](#)

Stable URL:

<http://0-www.jstor.org.skyline.ucdenver.edu/stable/2729411>

Page Count: 18



The problem, and our main point, is that a difference can be permanent...without being "significant" in other senses, such as for science or policy. And a difference can be significant for science or policy and yet be insignificant statistically, ignored by the less thoughtful researchers.

--McCloskey and Ziliak (1996, p. 97)

We here examine the alarming hypothesis that ordinary usage in economics takes statistical significance to be the same as economic significance.

--McCloskey and Ziliak (1996, p. 98)

Statistical computing software routinely provide t-statistics for every estimated coefficient. But that programs provide it does not mean that the information is relevant for science. We suspect that referees enforce the proliferation of meaningless t- and F-statistics, out of the belief that statistical and substantive significance are the same.

--McCloskey and Ziliak (1996, p. 102)



The elementary point that "[t]here is no sharp border between 'significant' and 'insignificant,' only increasingly strong evidence as the P-value decreases" (David S. Moore and George P. McCabe 1993, p. 473) is not found in most of the earlier books from which most economists learned statistics and econometrics. The old classic by W. Allen Wallis and Harry V. Roberts, *Statistics: A New Approach*, first published in 1956, is an exception:

It is essential not to confuse the statistical usage of "significant" with the everyday usage. In everyday usage, "significant" means "of practical importance," or simply "important." In statistical usage, "significant" means "signifying a characteristic of the population from which the sample is drawn," regardless of whether the characteristic is important. (Wallis and Roberts [1956] 1965, p. 385)



Magnitude and Significance

**OLS estimates of the relationship between years of schooling
and ln(wages)**

	Insignificant	Significant
BIG	.15 (.131)	.15 (.033)
Small	.0012 (.003)	.0012 (.00001)

Note: Standard errors are reported in parentheses.



OLS estimates of the relationship between years of schooling and ln(wages)

	Insignificant	Significant
BIG	.15 (.131)	.15 (.033)
Small	.0012 (.0033)	.0012 (.00001)

Note: Standard errors are reported in parentheses.

Why would we characterize these estimates as “BIG”? Big relative to what?



OLS estimates of the relationship between years of schooling and ln(wages)

	Insignificant	Significant
BIG	.15 (.131)	.15 (.033)
Small	.0012 (.0033)	.0012 (.00001)

Note: Standard errors are reported in parentheses.

Why would we characterize these estimates as “BIG”? Big relative to what?

1. Relative to previous estimates of the effect of schooling on wages...Cite the most important studies in this area (e.g., Angrist and Krueger 1991) or refer to a “stylized fact”.



WIKIPEDIA
The Free Encyclopedia

Stylized fact

From Wikipedia, the free encyclopedia

(Redirected from [Stylized facts](#))

Jump to: [navigation](#), [search](#)

In [social sciences](#), especially [economics](#), a **stylized fact** is a simplified presentation of an empirical finding.^[1] A stylized fact is often a broad generalization that summarizes some complicated statistical calculations, which although essentially true may have inaccuracies in the detail.



OLS estimates of the relationship between years of schooling and ln(wages)

	Insignificant	Significant
BIG	.15 (.131)	.15 (.033)
Small	.0012 (.0033)	.0012 (.00001)

Note: Standard errors are reported in parentheses.

Why would we characterize these estimates as “BIG”? Big relative to what?

2. Relative to the effect of other policies or to an established benchmark/goal. For instance, you could note that, according to your estimate, an additional year of schooling for Blacks would reduce the Black-White wage gap by 4 percentage points.



OLS estimates of the relationship between years of schooling and ln(wages)

	Insignificant	Significant
BIG	.15 (.131)	.15 (.033)
Small	.0012 (.0033)	.0012 (.00001)

Note: Standard errors are reported in parentheses.

Why would we characterize these estimates as “BIG”? Big relative to what?

3. Relative to the mean wage in your sample, or to the change in wages from 1970 to 2015. You’re looking for some way of anchoring the estimate and putting its magnitude into context...



OLS estimates of the relationship between years of schooling and ln(wages)

	Insignificant	Significant
BIG	.15 (.131)	.15 (.033)
Small	.0012 (.0033)	.0012 (.00001)

Note: Standard errors are reported in parentheses.

Why would we characterize these estimates as “^{small}”? Small relative to what?

The above strategies still apply. Often a good idea to point reader to the upper bound of the 95% CI. Even the upper bound of the 95% CI is small relative to...

Sometimes you just want to confirm that your estimates are reasonable...



Using data on 19- through 22-year-olds and an RD design, Carpenter and Dobkin (2009) found that reaching the minimum legal drinking age was associated with a 21% increase in alcohol consumption and a 15% increase in traffic fatalities

The implied elasticity from these estimates is 0.71 (i.e., 0.15/0.21)

Restricting our sample to 19- through 22-year-olds, we (Anderson et al. 2013) found that legalizing medical marijuana was associated with a 15% decrease in alcohol consumption and a 12% decrease in traffic fatalities

The implied elasticity from these estimates is 0.80 (i.e., 0.12/0.15)

References

Anderson, D. Mark, Benjamin Hansen, and Daniel I. Rees. 2013. "Medical Marijuana Laws, Traffic Fatalities, and Alcohol Consumption." *Journal of Law and Economics*, 56 (2): 333-369.

Carpenter, Christopher and Carlos Dobkin. 2009. "The Effect of Alcohol Consumption on Mortality: Regression Discontinuity Evidence from the Minimum Drinking Age." *AEJ: Applied*, 1(1): 164-182.



Data Mining

Much of this lecture is based on Angrist and Pischke (2010). It also refers to Ehrlich (1975, 1977), Leamer (1983), and other researchers.

References

Angrist, Joshua D. and Jorn-Steffen Pischke. 2010. "The Credibility Revolution in Empirical Economics: How Better Research Design Is Taking the Con out of Econometrics." *Journal of Economic Perspectives*, Vol. 24(2), pp. 3-30.

Donohue, John and Justin Wolfers. 2005. "Uses and Abuses of Statistical Evidence in the Death Penalty Debate", *Stanford Law Review*, Vol. 58, pp. 787-

Donohue, John and Justin J. Wolfers. 2006. "The Death Penalty: No Evidence For Deterrence." *The Economists' Voice*, Vol. 3 (5), Article 3.

Ehrlich, Isaac. 1975. "The Deterrent Effect of Capital Punishment: A Question of Life and Death." *The American Economic Review*, Vol. 65 (3), pp. 397-417.

Ehrlich, Isaac. 1977. "Capital Punishment and Deterrence: Some Further Thoughts and Additional Evidence." *Journal of Political Economy*, Vol. 85 (4), pp. 741-788

Leamer, Edward E. 1983. "Let's Take the Con Out of Econometrics." *The American Economic Review*, Vol. 73 (1), pp. 31-43.

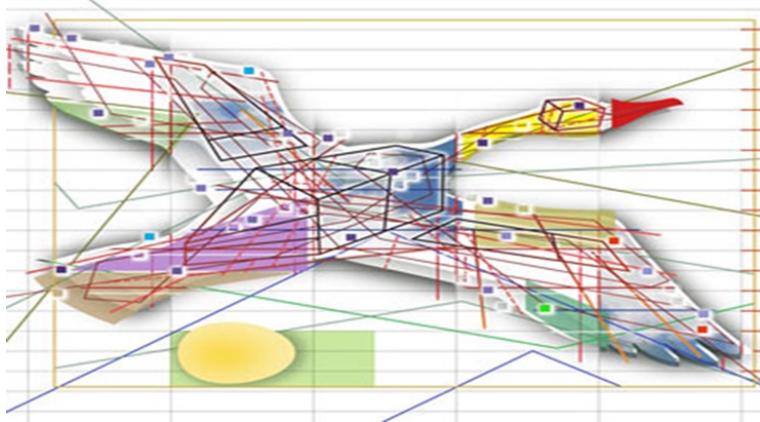
Mocan, H. Naci and R. Kaj Gittings. 2003. "Getting Off Death Row: Commuted Sentences and the Deterrent Effect of Capital Punishment", *Journal of Law and Economics*, Vol. 46(2), pp. 453-478.

Data Mining



If you torture the data long enough, nature will confess.

--Ronald Coase



This is a sad and decidedly unscientific state of affairs we find ourselves in. Hardly anyone takes data analyses seriously. Or perhaps more accurately, hardly anyone takes anyone else's data analyses seriously. Like elaborately plumed birds who have long since lost the ability to procreate but not the desire, we preen and strut and display our t-values.

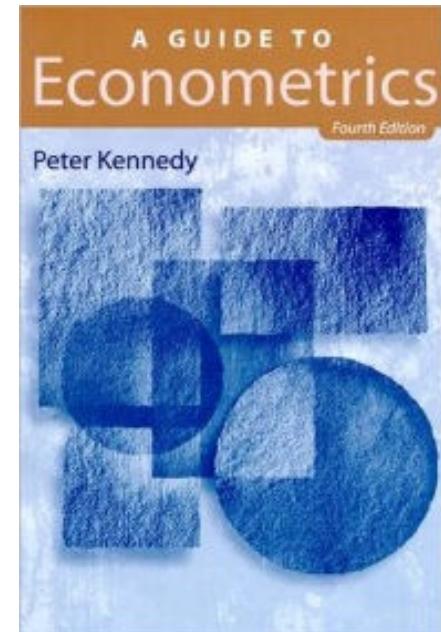
--Edward Leamer (1983, p. 37)



“Classic” data mining or p-hacking

The terminology “data mining” is often used in the context of pre-test bias. In particular, researchers often run a large number of different regressions on a body of data looking for significant t statistics (at ,say, the 5% level). Using this approach invalidates traditional hypothesis-testing procedures because such data mining is likely by chance to uncover significant t statistics; i.e., the final results chosen are much more likely to embody Type I error than the claimed 5%.

--From Peter Kennedy's *Guide to Econometrics* (4th edition)





Three signs you may be data mining

1. You only report estimates of the following equation:

$$y_{it} = \alpha + \beta_1 X_{1it} + \beta_2 X_{2it} + \dots \beta_k X_{kit} + \varepsilon_{it} .$$

However, estimating this equation produces very different results:

$$\ln(y_{it}) = \alpha + \beta_1 X_{1it} + \beta_2 X_{2it} + \dots \beta_k X_{kit} + \varepsilon_{it} .$$



2. You only report estimates of the following equation:

$$y_{it} = \alpha_0 + \alpha_1 Z_{it-1} + \beta_1 X_{1it} + \beta_2 X_{2it} + \dots \beta_k X_{kit} + \varepsilon_{it}.$$

However, estimating any of these equations produces very different results:

$$y_{it} = \alpha_0 + \alpha_1 Z_{it} + \beta_1 X_{1it} + \beta_2 X_{2it} + \dots \beta_k X_{kit} + \varepsilon_{it}$$

$$y_{it} = \alpha_0 + \alpha_1 Z_{it-2} + \beta_1 X_{1it} + \beta_2 X_{2it} + \dots \beta_k X_{kit} + \varepsilon_{it}$$

$$y_{it} = \alpha_0 + \alpha_1 Z_{it-1} + \alpha_2 Z_{it-2} + \beta_1 X_{1it} + \beta_2 X_{2it} + \dots \beta_k X_{kit} + \varepsilon_{it}$$

$$y_{it} = \alpha_0 + \alpha_1 Z_{it} + \alpha_2 Z_{it-1} + \alpha_3 Z_{it-2} + \beta_1 X_{1it} + \beta_2 X_{2it} + \dots \beta_k X_{kit} + \varepsilon_{it}.$$



3. You only report estimates of the following equation for respondents ages 12-15:

$$y_{it} = \alpha_0 + \alpha_1 Z_{it-1} + \beta_1 X_{1it} + \beta_2 X_{2it} + \dots \beta_k X_{kit} + \varepsilon_{it}.$$

However, you get very different results if your sample includes 16-year-olds.

A few more behaviors that can be described as data mining:

- Stop project because of null results
 - Decide against submitting to a journal
- Modifying your original hypothesis to better match results



Data mining is not the same thing as doing robustness checks

Table 5. Robustness Checks

	Analgesics	Hallucinogens
1) Any disorder (baseline)	0.039** (0.0161)	0.0217** (0.0092)
2) Mild disorder	0.0162 (0.0137)	0.0106** (0.0051)
3) Moderate disorder	0.0095 (0.0059)	0.0033 (0.0043)
4) Severe disorder	0.0122** (0.006)	0.0086*** (0.0023)
5) Any disorder, no imputed outcomes	0.0390** (0.0168)	0.0220** (0.0093)
6) Any disorder, conditional on past year use	0.8232** (0.3548)	0.5748* (0.3056)
7) Any disorder, linear probability model	0.0386** (0.0170)	0.0202** (0.0095)
8) Any disorder, probit	0.0391** (0.0161)	0.0184** (0.0078)
9) Any disorder, economic conditions measured by employment rate	-0.0120 (0.0150)	-0.0172** (0.0080)
10) Any disorder, economic conditions measured by log GDP per capita	-0.00882* (0.00504)	-0.00232 (0.00156)
11) Any disorder, control for lagged unemployment	0.0385** (0.0163)	0.0214** (0.0090)



Among non-economists, data mining is not a dirty secret.

Big Data: Exploratory Data Mining in Behavioral Research

Arizona State University

Tempe, Ariz.

June 1-5, 2015

Monday-Friday



This ATI provides an overview of recent methodological advances in exploratory data mining for the analysis of psychological and behavioral data. In contrast to traditional hypothesis-driven approaches to analysis, exploratory data mining enables investigators to assess the predictive value of all possible combinations of variables in a data set. Data mining has emerged in recent years as a major area of statistical research and practice and is increasingly employed by psychologists and other behavioral scientists.

Exploratory data mining techniques are particularly useful for the analysis of very large data sets, as can arise in clinical, survey, psychometric and genomic research. These techniques are often a natural follow-up to standard multivariate analyses in cases in which investigators have either: (1) obtained significant results and seek to know whether there are other important patterns in the data; (2) obtained no notable results and wonder whether there are any important patterns to be found; or (3) developed questions that are too general or imprecisely formulated to be addressed through hypothesis testing.

The five-day course will cover the conceptual bases and strategies of exploratory data mining, and will review leading current techniques and software, including those based on canonical regression models and on pattern searches with recursive partitioning (see a [tentative course schedule](#)). Participants should also bring their own laptop computer equipped with SPSS or SAS and with R (available for free [online](#)). The primary software for instruction will be SAS and R, although code will also be provided for SPSS and instructors will be familiar with that and other programs.

The course will be directed by John McArdle (University of Southern California). Other instructors are expected to include: Gilbert Ritschard (University of Geneva, Switzerland), George Marcoulides (University of California, Riverside) and Kevin Grimm (Arizona State University).



Capital Punishment: A Cautionary Tale

Note that, in theory, the relationship between capital punishment and violent crime could be negative. On the other hand, would-be criminals could be completely oblivious to the potential consequences of their actions.

Ehrlich (1975, 1977) estimated the relationship between violent crime and executions. He concluded that capital punishment acted as a deterrent. Specifically, he wrote, “[c]ontrary to strong inferences of many death-penalty researchers who interpreted their evidence as a categorical denial of the deterrence hypothesis, this study as well as a related analysis based on time-series data develop evidence not inconsistent with that hypothesis.”

Ehrlich's work was *very* influential. It was cited repeatedly in *Gregg v. Georgia* (United States Supreme Court, 1976), which reaffirmed the use of the death penalty in the United States. It also spawned a large number of follow-up studies.



TABLE 2—VARIABLES USED IN THE REGRESSION ANALYSIS, ANNUAL OBSERVATIONS 1933–69

Variable		Mean (Natural Logarithms)	Standard Deviation	Arithmetic Mean
y_1	{ $(Q/N)^0$ = Crime rate: offenses known per 1,000 civilian population.	-2.857	0.156	0.058
y_1	P^0a = Probability of arrest: percent of offenses cleared.	4.997	0.038	89.835
y_1	$P^0c a$ = Conditional probability of conviction: percent of those charged who were convicted of murder. ^a	3.741	0.175	42.733
y_1	$P^0e c$ = Conditional probability of execution; PXQ_1 = the number of executions for murder in the year $t+1$ as a percent of the total number of convictions in year t . ^b	0.176	1.749	2.590
X_1	L = Labor force participation: fraction of the civilian population in the labor force.	-0.546	0.030	0.579
X_1	U = Unemployment rate: percent of the civilian labor force unemployed.	1.743	0.728	7.532
X_1	A = Fraction of residential population in the age group 14–24.	-1.740	0.118	0.177
X_1	Y_p = Friedman's estimate of (real) permanent income per capita in dollars.	6.868	0.338	1012.35
X_1	T = Chronological time (years): 31–37.	2.685	0.867	19.00
X_2	NW = Fraction of nonwhites in residential population.	-2.212	0.063	0.110
X_2	N = Civilian population in 1,000s.	11.944	0.161	155,853
X_2	$XGOV$ = Per capita (real) expenditures (excluding national defense) of all governments in million dollars.	-7.661	0.501	.000532
X_2	$XPOL_{-1}$ = Per capita (real) expenditures on police in dollars lagged one year. ^a	2.114	0.306	8.638

^a The figures for $P^0c|a$ (1933–35) and $XPOL$ (all the odd years 1933–51) were interpolated via an auxiliary regression analysis.

^b The actual number of executions 1968, 1969, and 1970 was zero. However, the numbers were assumed equal to 1 in each of these years in constructing the value of PXQ_1 in 1967–69.

Endogenous Variables

Instruments



The New York Times

Does Death Penalty Save Lives? A New Debate

By [ADAM LIPTAK](#)

Published: November 18, 2007



For the first time in a generation, the question of whether the death penalty deters murders has captured the attention of scholars in law and economics, setting off an intense new debate about one of the central justifications for capital punishment.

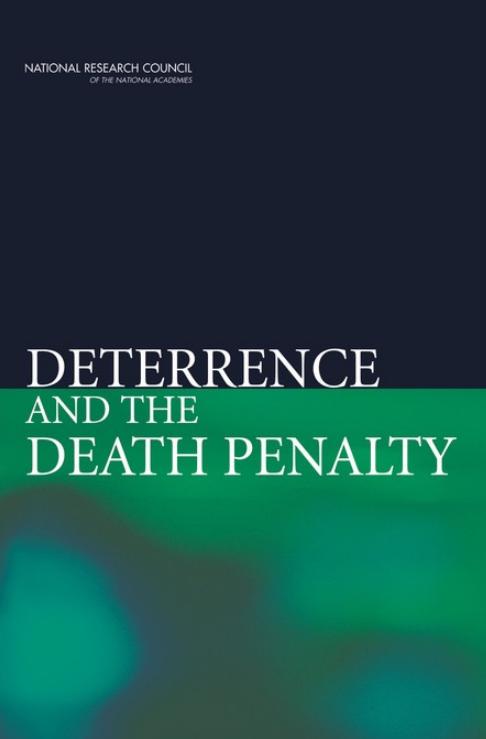
According to roughly a dozen recent studies, executions save lives. For each inmate put to death, the studies say, 3 to 18 murders are prevented.

The studies, performed by economists in the past decade, compare the number of executions in different jurisdictions with homicide rates over time — while trying to eliminate the effects of crime rates, conviction rates and other factors — and say that murder rates tend to fall as executions rise. One influential study looked at 3,054 counties over two decades.

“I personally am opposed to the death penalty,” said H. Naci Mocan, an economist at Louisiana State University and an author of a study finding that each execution saves five lives. “But my research shows that there is a deterrent effect.”



An expert panel convened by the National Academy of Sciences concluded that existing research “is not informative about whether capital punishment decreases, increases or has no effect on homicide rates,” and that such studies “should not influence policy judgments about capital punishment.”





Ronnie Hillman fits in Broncos offense, moves up to No. 2 on depth chart



Navajo Nation not lifting San Juan River closure after EPA OK's water



Police: Louisiana trooper shot in head



Home

Colorado Breaking News, Sports, Weather, Traffic, Jobs

Story

COLORADO BREAKING NEWS, SPORTS, WEATHER, TRAFFIC, JOBS

No credible evidence on whether death penalty deters, experts say

By Michael Booth

The Denver Post

POSTED: 06/03/2013 12:01:00 AM MDT | UPDATED: 2 YEARS AGO

248 COMMENTS

Go ahead and stake out your opinion on whether a state's death-penalty law will deter future murderers.

But don't pretend that opinion is based on any remotely credible evidence, according to a consensus of criminologists, economists and other academics who have reviewed deterrence studies from both sides and officially declared them useless.

The National Academy of Sciences picked apart decades of deterrence research last year and recommended "that these studies not be used to inform deliberations" on capital punishment. In a blunt report, the academy's National Research Council noted it had made a similar survey 30 years before and was "disappointed" to learn each study since was equally futile.

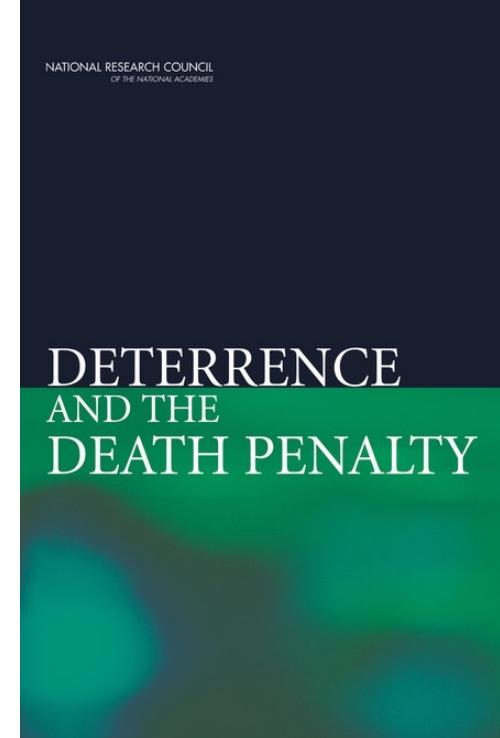


Nathan Dunlap on May 1, 2013.

Sep 11:



An expert panel convened by the National Academy of Sciences concluded that concluded that existing research “is not informative about whether capital punishment decreases, increases or has no effect on homicide rates,” and that such studies “should not influence policy judgments about capital punishment.”



What happened to Mocan and Gittings (2003)? Donohue and Wolfers (2005) re-examined their data...



TABLE 6: ESTIMATING THE IMPACT OF EXECUTIONS ON MURDER RATES:
RE-ANALYZING MOCAN AND GITTINGS: 1977-1997

	Dependent Variable:						
	<i>Annual Homicides per 100,000 Residents_{s,t}</i>			<i>Log Homicide Rate_{s,t}</i>			
	(1)	(2)	(3)	(4)	(5)	(6)	(7)
Panel A: Mocan and Gittings Results: Replication							
Executions _{t-1} per Death Sentence _{t-7}	-0.60 [*] (.35)			-0.63 [*] (.34)	-0.63 ^{**} (.29)	-0.05 [*] (.03)	-0.05 [*] (.03)
Pardons _{t-1} per Death Sentence _{t-7}		0.69 ^{**} (.32)		0.73 ^{**} (.30)		0.11 ^{***} (.03)	
Death Row Removals _{t-1} per Death Sentence _{t-6}			0.17 ^{**} (.07)		0.18 ^{**} (.07)		0.02 ^{**} (.01)
Sample (1984-1997)	680	693	695	679	690	679	690
Panel B: Correcting Programming Errors							
Executions _{t-1} per Death Sentence _{t-7}	-0.50 (.34)			-0.52 (.33)	-0.59 (0.39)	-0.01 (0.03)	-0.02 (0.02)
Pardons _{t-1} per Death Sentence _{t-7}		0.63 [*] (.34)		0.71 ^{**} (.30)		0.09 ^{***} (0.03)	
Death Row Removals _{t-1} per Death Sentence _{t-6}			0.24 ^{***} (.08)		0.17 [*] (0.09)		0.01 (0.01)
Sample (1984-1997)	679	692	691	677	636	677	636
Panel C: Measuring Deterrence Variables with a One-Year Lag on Full Sample							
Executions _{t-1} per Death Sentence _{t-7}	0.03 (0.14)			0.01 (0.13)	0.01 (0.14)	0.01 (0.01)	0.01 (0.01)
Pardons _{t-1} per Death Sentence _{t-1}		0.41 ^{***} (.13)		0.41 ^{***} (0.13)		0.05 ^{***} (0.01)	
Death Row Removals _{t-1} per Death Sentence _{t-1}			0.02 (0.03)		0.02 (0.03)		0.002 (.002)
Sample (1978-1997)	986	984	921	977	918	977	918

t - 7

t - 1



given the paucity of evidence on how these expectations are formed, there seems little reason to strongly prefer one specification over the other. Thus, in Panel C, we rerun their regressions but note Zimmerman's argument that "any truly meaningful (subjective) assessment a potential murderer makes . . . is likely to be based upon the most recent information available to him/her."⁶⁶

an artificial negative correlation between execution and homicide rates.

64. Mocan & Gittings, *supra* note 11, at 478. While their data runs from 1977 to 1997, their complicated lag structure means that they can only estimate effects from 1984 onward.

65. Two types of coding errors were discovered. First, the authors attempted to drop all observations where the explanatory variable was the ratio of a positive value to zero but ended up both dropping the prior observation and including the variable they intended to drop, coded as the ratio of the numerator to 0.99. Second, in Models 3, 5, and 6, the execution rate was defined relative to the number of death sentences six years prior instead of seven years prior, as they did in their other specifications (and described in their text).|

66. Zimmerman, *supra* note 11, at 170.



UNPACKING P-HACKING AND PUBLICATION BIAS*

Abel Brodeur^a, Scott Carrell^{b†}, David Figlio^c, Lester Lusher^d

^aUniversity of Ottawa and IZA

^bUniversity of California, Davis, NBER, and IZA

^cNorthwestern University, NBER, and IZA

^dUniversity of Hawaii at Manoa and IZA

This draft: April 25, 2021

These authors examined data from 396 manuscripts submitted to the JHR during the period 2013-2018. Ninety-eight were desk rejected, 110 were rejected based on reports, and 94 were published.



Their Motivation

Studies have illustrated how the distribution of test statistics from published manuscripts lump at certain significance thresholds. Little is known about the underlying mechanisms driving these findings:

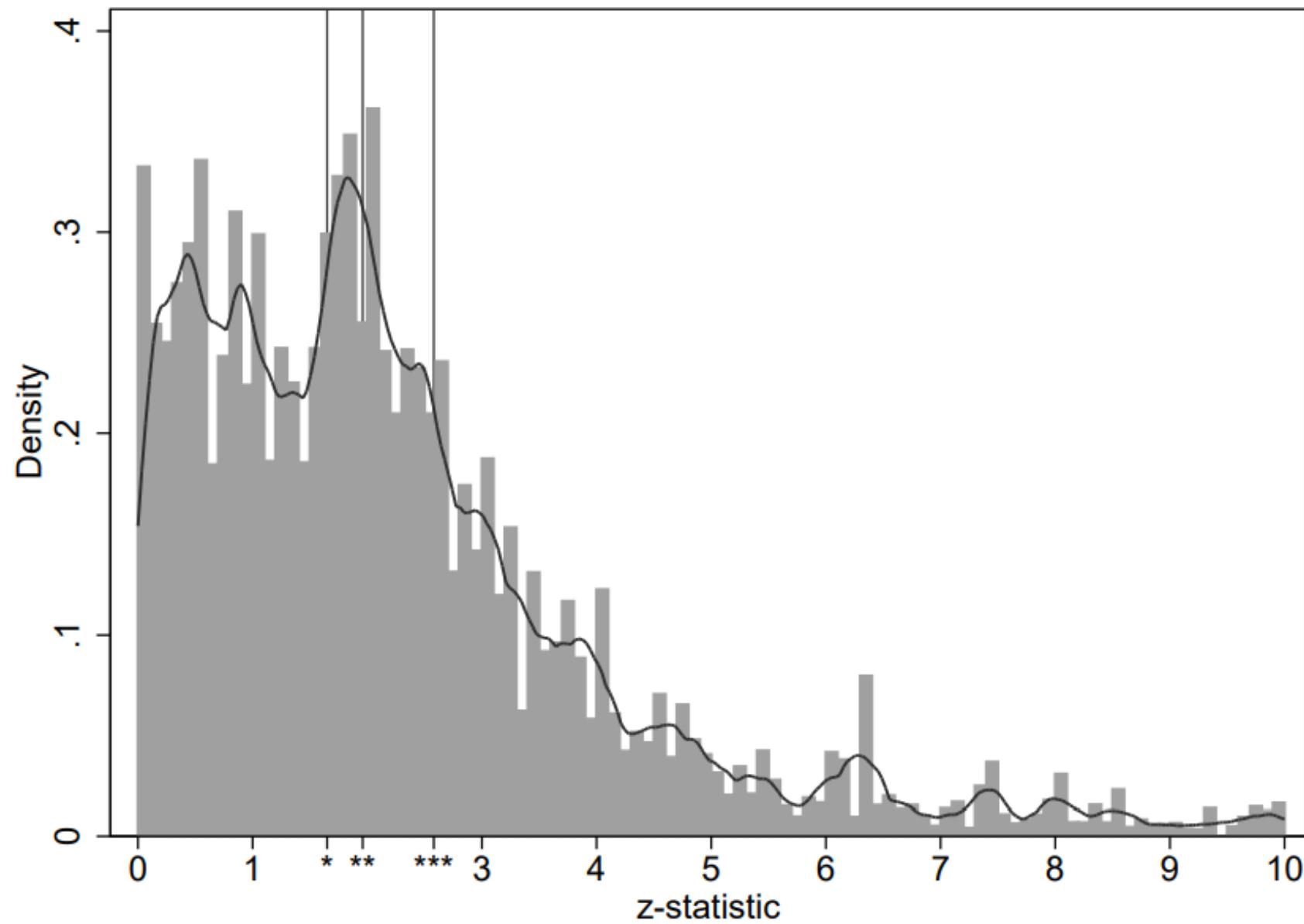
- (1) authors may engage in certain behavior to attain “desirable” p-values (p-hacking)
- (2) and/or the peer review process may favor statistical significance (publication bias).

This study is the first to use data from journal submissions to identify these mechanisms.



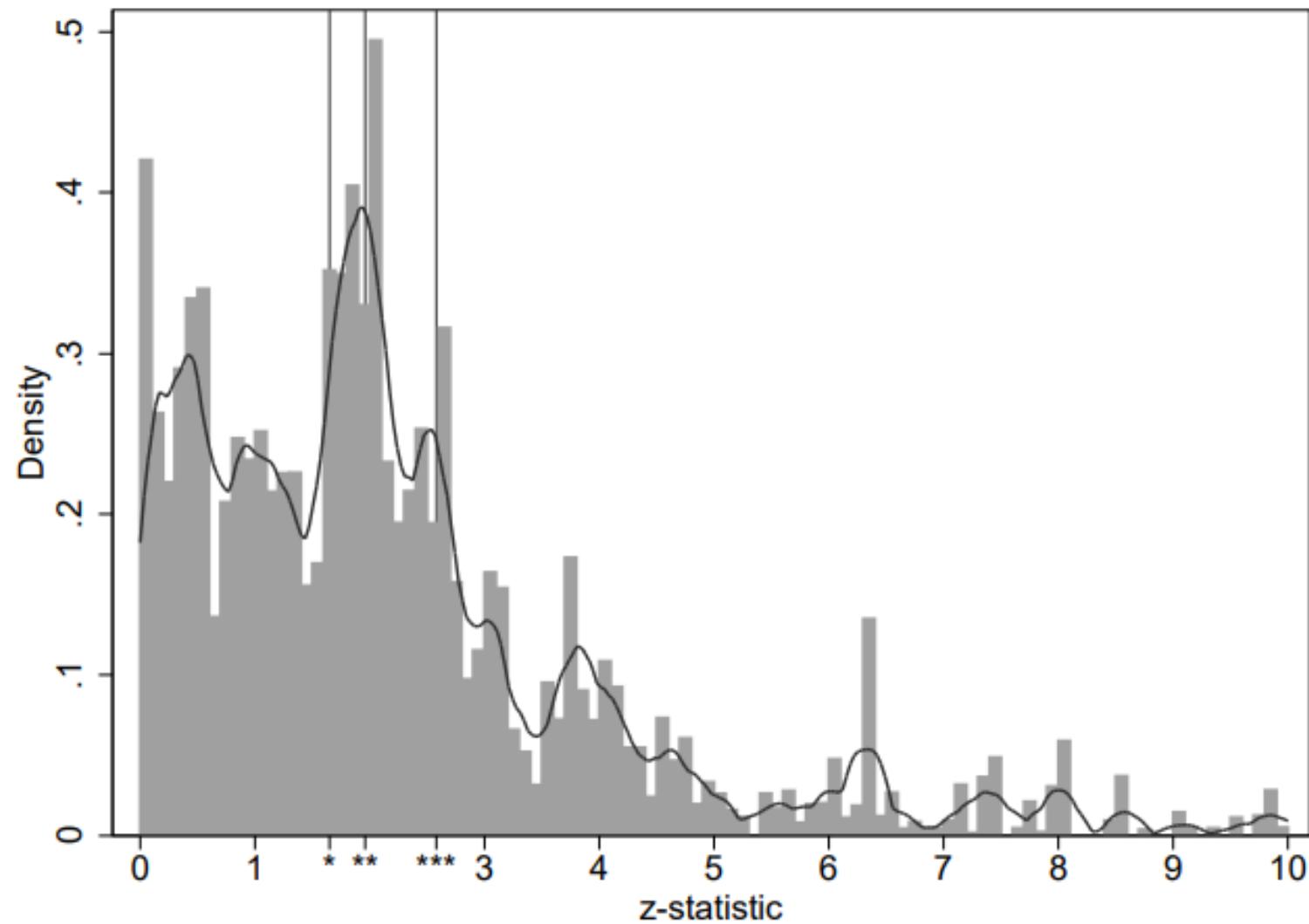
Their Findings

- We first find that initial submissions display significant bunching [at well-known thresholds of statistical significance], suggesting prior findings cannot be strictly attributed to a bias in the peer review process
- Desk rejected manuscripts display greater heaping than those sent for review, suggesting editors on average “sniff out” marginally significant results
- Reviewer recommendations, on the other hand, are swayed significantly by statistical thresholds...





(a) Desk rejections





Q&A (\approx 5 minutes)

Doing Applied Research

MIXTAPE TRACK



**Section 2. Practical tips for writing your applied paper:
Sensitivity and heterogeneity analyses**

Sensitivity Analysis



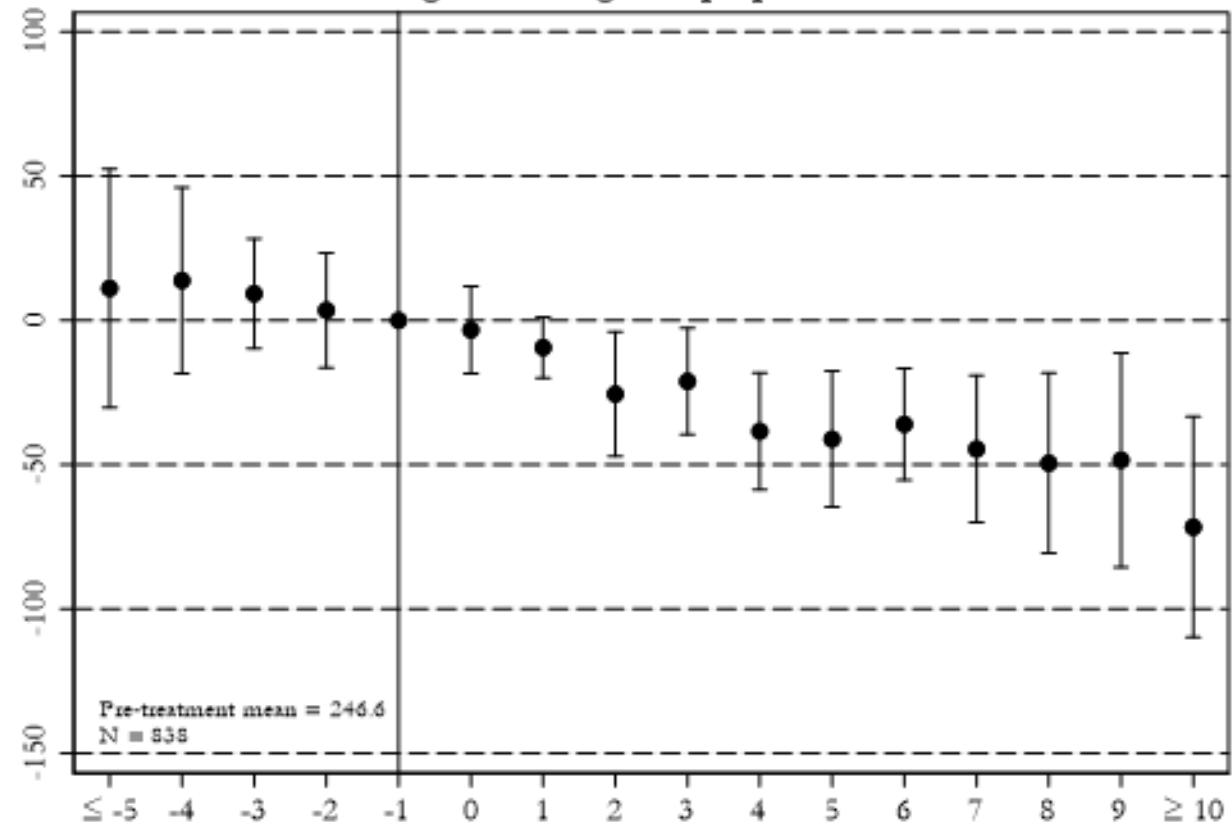
- Why do a sensitivity analysis?
 - Reassure the reader that you're not p-hacking
 - Good empirical work is transparent and easily replicated
 - Provide a range of estimates
 - Unless one particular specification is preferred, readers will expect a range of estimates
 - Often, theory does not dictate specification
 - Facilitate comparison of your estimates to those reported by previous researchers
 - To weight or not to weight?
 - Solon et al. (2015)
 - Unit-specific linear trends
 - TWFEs with staggered treatment

Reference

Solon, Gary, Steven J. Haider, and Jeffrey M. Wooldridge. 2015. "What Are We Weighting For?" *JHR*, 50(2): 301-316.

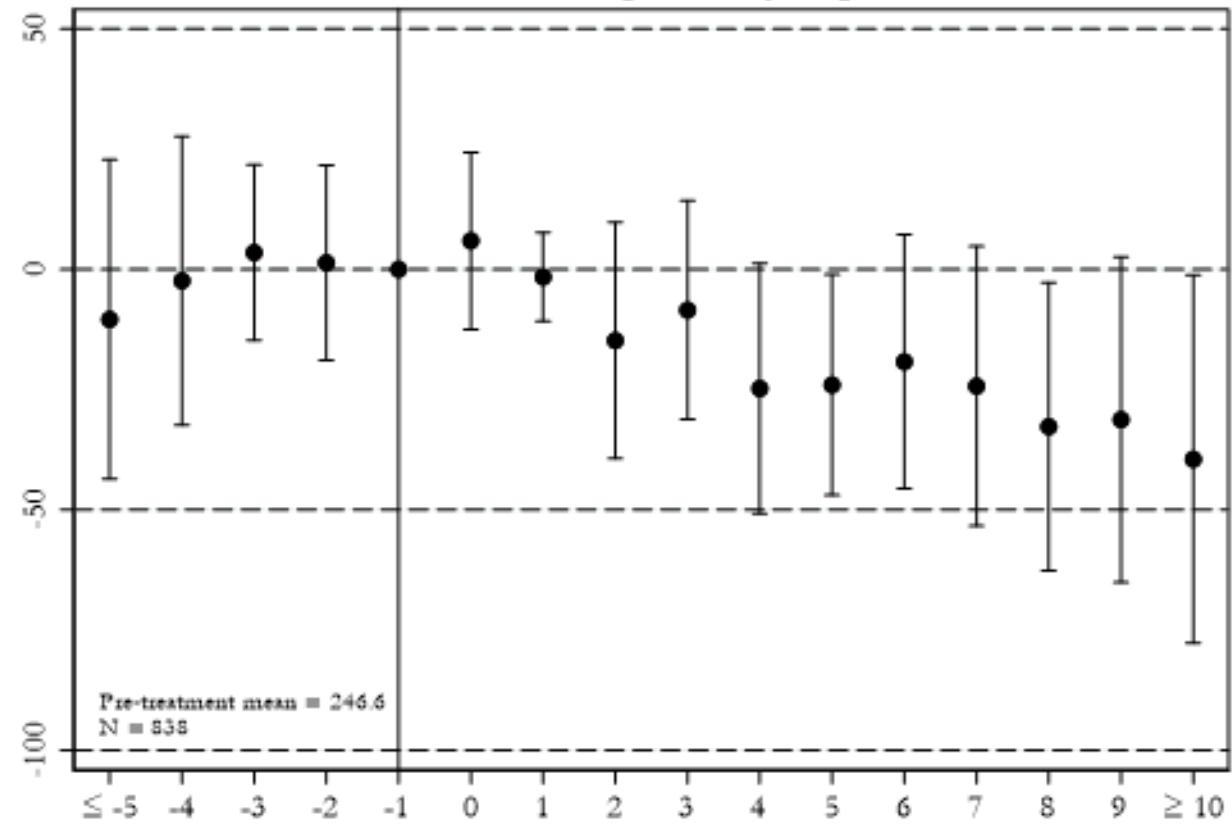


Panel A. Controlling for Timing Group-Specific Linear Time Trends



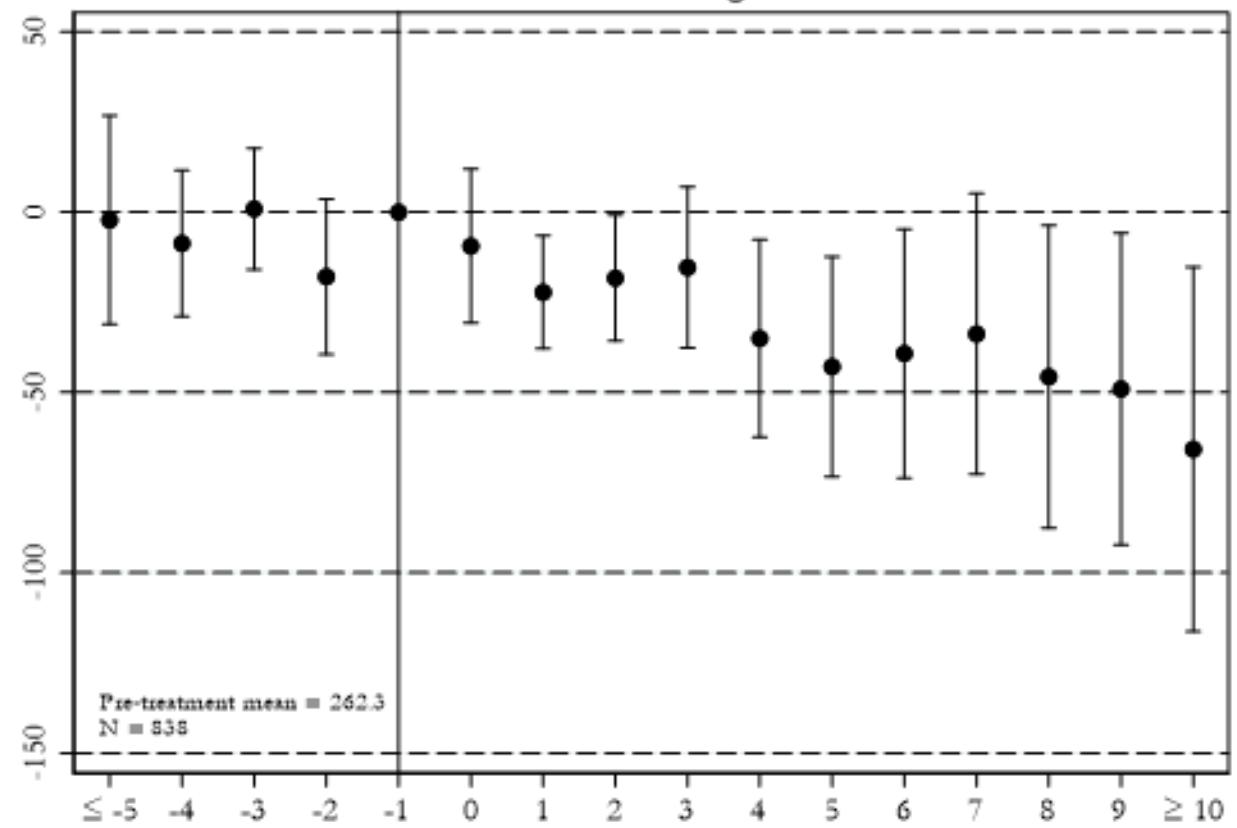


Panel B. Controlling for Dairy Inspections





Panel C. Unweighted





Panel D. Dropping Years After 1899

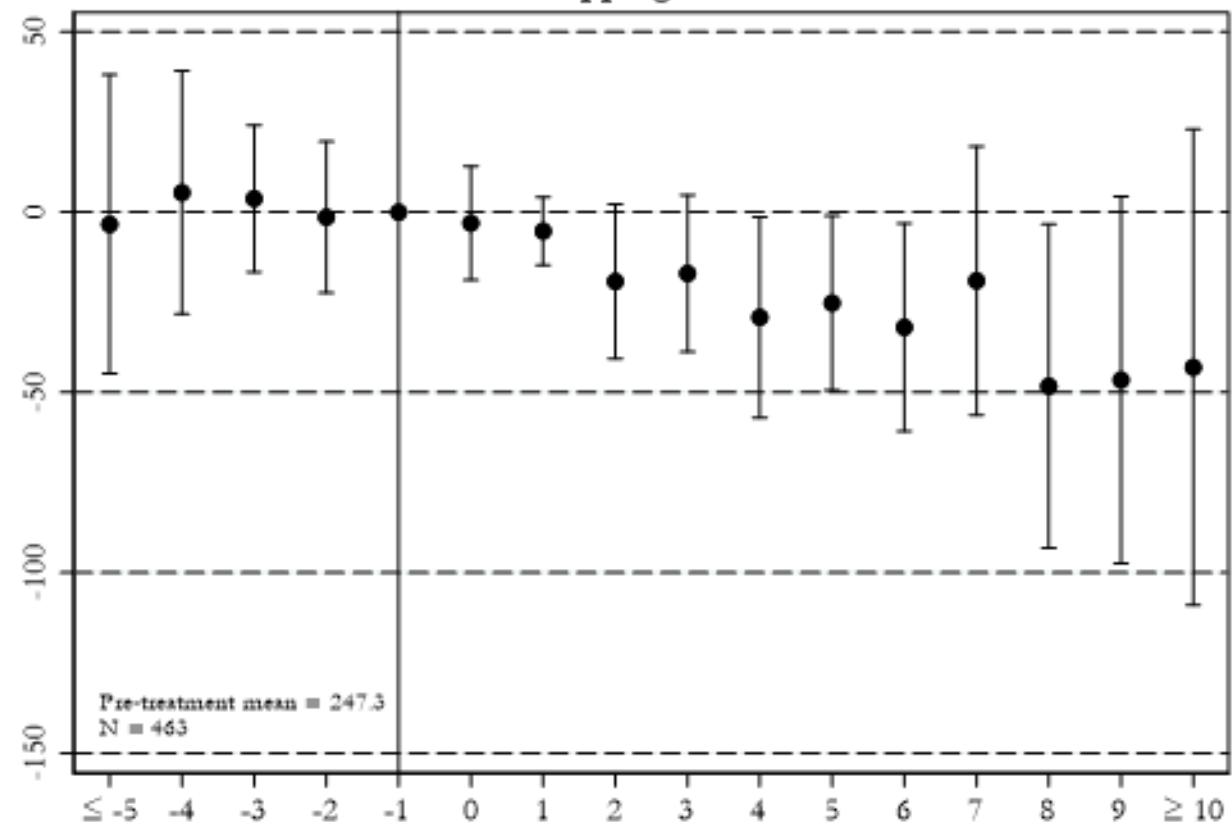
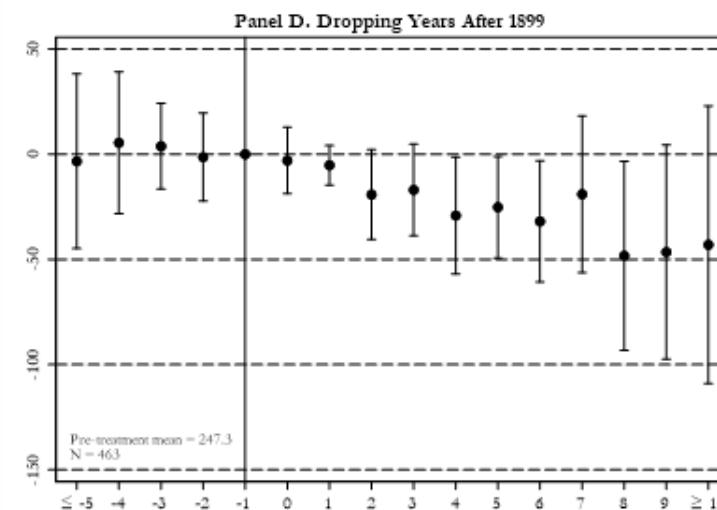
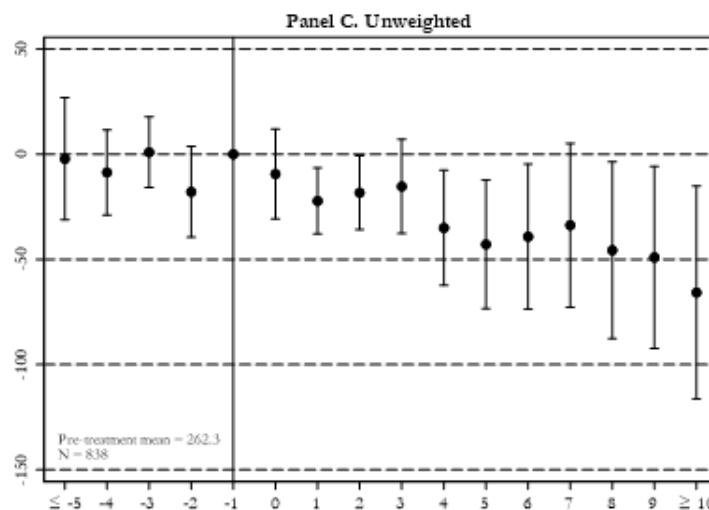
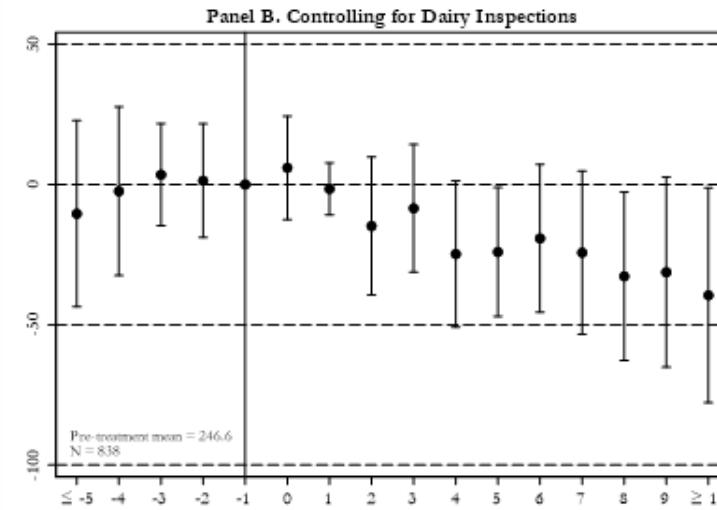
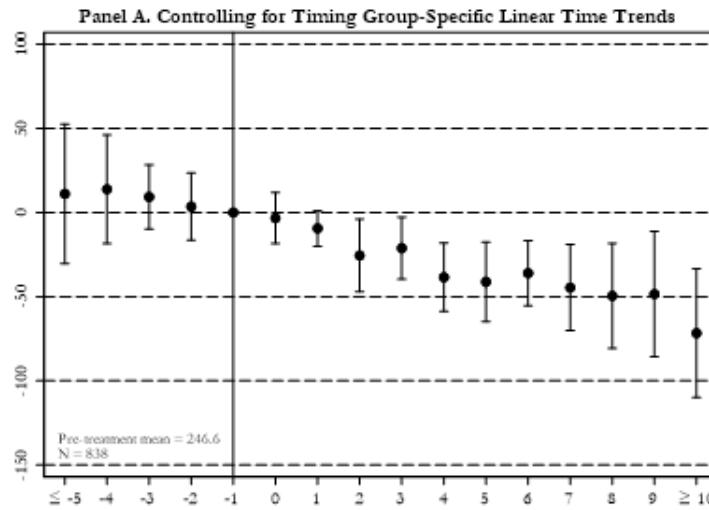




Figure 4. Milk Inspections and Waterborne Mortality: Sensitivity Analysis



The goal is to be thorough without getting bogged down. It is easy to get bogged down...



A few more thoughts about falsification tests or adopting an alternative ID strategy

- Goal
 - The goal is to convince the reader that your estimate is causal
- Describe the threat to identification
 - Describing a variable as “potentially endogenous” is not enough
 - Describing the threat to identification sets up the falsification test or alternative ID strategy



Q&A (\approx 5 minutes)

Doing Applied Research

MIXTAPE TRACK



**Section 2. Practical tips for writing your applied paper:
How to write an effective conclusion**



The Conclusion

- Repeat yourself, but don't repeat yourself!
 - The conclusion should contain many of the same elements found in the introduction.
 - But, it should be shorter and contain fewer background details.
 - (i) Open by again providing a general motivation
 - If you used some eye-catching stats as motivation in your intro, try to find some different stats for your conclusion that provide the same motivation
 - Consider the following example from Anderson and Sabia (2018, *J of Law and Economics*) who estimated the relationship between safe-storage gun laws and school shootings in the United States:
 - First paragraph of introduction: "School shootings, such as the recent high-profile events in Chardon, Ohio; ...and Santa Fe, Texas are usually committed by students under the age of 18. In fact, between 2012 and 2015, approximately 70 percent of shootings in K-12 schools in the United States were committed by minors."
 - First paragraph of conclusion: "The National Poll on Children's Health recently indicated that 1 in 4 parents are 'very concerned' about school violence for their children. These fears are perhaps driven by the fact that school shootings have been reported at nearly a weekly rate since 2012."
 - (ii) Reiterate the specific motivation for scholars working in this area. Remind readers of the contribution you are making to the literature.
 - (iii) Briefly describe your setting, data, and empirical/identification strategy but you do not have to go into as much detail as earlier
 - (vi) Results discussion.
 - This is the one place in the conclusion where I will go into more depth than in the introduction
 - Discuss your principal findings but do not be afraid to also remind the reader of interesting sub-analyses (e.g., heterogeneous effects, particularly convincing robustness checks, supplementary data analyses, etc.)



The Conclusion

- The last paragraph of the conclusion is *really* important. Finish with a bang!
- I find this to be one of the more difficult paragraphs to write in a paper. It can feel daunting, but you have options...
 - (i) Opine about future research
 - (ii) Highlight a contribution that didn't make it into your introduction
 - Historical U.S. studies that relate estimates to developing countries today
 - "Infant Health and Later-Life Labor Market Outcomes," Lazuka (2018, *J of Human Resources*)
 - Milk inspections paper cited earlier (Anderson et al. 2022, NBER WP No. 30063)
 - (iii) Policy relevance
 - "The Impact of Mass Shootings on Gun Policy," Luca et al. (2020, *J of Public Economics*)
 - "Accidents Will Happen? Unintentional Childhood Injuries and the Effects of Child Care Regulations," Currie and Hotz (2004, *J of Health Economics*)
 - (iv) List caveats to your study
 - Not as big a fan of this for the final paragraph. I like including caveats in the conclusion, but maybe not what you want to leave readers with as your final word.
 - That said, this is a common approach for the concluding paragraph and here is a nice example:
 - "Is Legal Pot Crippling Mexican Drug Trafficking Organisations? The Effect of Medical Marijuana Laws on U.S. Crime," Gavrilova et al. (2017, *Economic Journal*)
 - Or some combination of (i)-(iv)



The Conclusion

- How long should it be?
 - Shorter than your introduction.
 - The last 15 papers I have written as full-length articles for economics journals have had conclusions of the following lengths (size 12 font, double-spaced):

1-2 pages	2-3 pages	3+ pages
12 papers	3 papers	0 papers



Q&A (\approx 10 minutes)