

# Econometrics 1: Tutorial IV

Heikki Korpela

December 4, 2023

# Contents

1. Practical matters
2. Issues with data
3. Instrumental variables

# 1. Practical matters

- ▶ An extra tutorial will be held on Friday (Dec 8) at 2.15pm
- ▶ My apologies if this is inconvenient for you
  - ▶ I had to find a time quickly that seemed like it would have the least clash with lectures for other core Economics courses
  - ▶ Next Monday you will have the last regular lecture with the professor
- ▶ You can also email me (as many of you have), and also schedule a meeting if you want to discuss any problems

## 2. Issues with data

- ▶ The dataset for problem 2 is actually quite instructive: it illuminates (depressingly) common features of real data.
- ▶ You usually need to read metadata rather carefully to learn what the variables actually mean.
- ▶ Technical issues and corrupting data by accident are quite common.
- ▶ The original `schooling.txt` had an issue, which could cause it to fail to import.
  - ▶ The first row of data had different numbers of whitespace between column names.
  - ▶ Software can easily get confused by things like this. There is no universal rule for how to cope with such inconsistencies, so the software often makes a best guess.

# How to cope?

- ▶ **Look at your data.** Calculate summary statistics like the mean, and check if they make rough sense.
  - ▶ Example: a variable holding values like 50 000 is never going to be in logs. (Base R's numeric value can hold values up to about  $\exp(709.78)$ .)
- ▶ If you run into trouble or suspect issues, check out options and help pages for the functions you use to import data.
- ▶ Read the metadata. (Generally not required for this course, the assignments should have the descriptions required or refer you to the slides/textbook.)

# The journey to metadata 1/6

*ILLUSTRATION: ESTIMATING THE RETURNS TO SCHOOLING*

159

In this section we use data on 3010 men taken from the US National Longitudinal Survey of Young Men, also employed in Card (1995). In this panel survey, a group of

The bibliography reference in the textbook description of the data.

# The journey to metadata 2/6

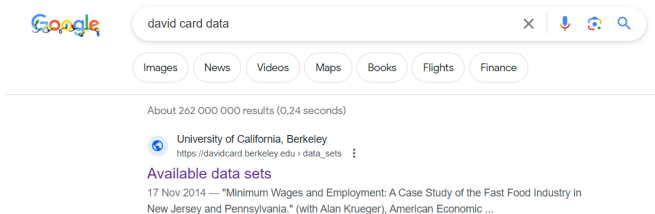
472

BIBLIOGRAPHY

- Cameron, A. C. and Windmeijer, F. A. G. (1996), *R*-squared Measures for Count Data Regression Models with Applications to Health Care Utilization, *Journal of Business and Economic Statistics*, 14, 209–220.
- Cameron, A. C. and Windmeijer, F. A. G. (1997), An *R*-squared Measure of Goodness of Fit for Some Common Nonlinear Regression Models, *Journal of Econometrics*, 77, 329–342.
- Campbell, J. Y. and Perron, P. (1991), Pitfalls and Opportunities: What Macroeconomists Should Know about Unit Roots. In: O. Blanchard and S. Fisher, eds, *NBER Macroeconomics Annual*, 6, 141–201, MIT Press, Cambridge.
- Campbell, J. Y. and Shiller, R. J. (1991), Yield Spreads and Interest Rate Movements: A Bird's Eye View, *Review of Economic Studies*, 58, 495–514.
- Campbell, J. Y. and Shiller, R. J. (1998), Valuation Ratios and the Long-Run Stock Market Outlook, *Journal of Portfolio Management*, 24, 11–26.
- Campbell, J. Y. and Thompson, S. B. (2008), Predicting Excess Stock Returns Out of Sample: Can Anything Beat the Historical Average?, *Review of Financial Studies*, 21, 1509–1531.
- Campbell, J. Y., Lo, A. W. and MacKinlay, A. C. (1997), *The Econometrics of Financial Markets*, Princeton University Press, Princeton, NJ.
- Canova, F. (1995), The Economics of VAR Models. In: K. D. Hoover, ed., *Macroeconometrics: Developments, Tensions and Prospects*, Kluwer Academic Publishers, Boston, MA, 57–97.
- Canova, F. (2007), *Methods for Applied Macroeconomic Research*, Princeton University Press, Princeton, NJ.
- Card, D. (1995), Using Geographical Variation in College Proximity to Estimate the Return to Schooling. In: L. N. Christofides, E. K. Grant and R. Swidinsky, eds, *Aspects of Labour Market Behaviour: Essays in Honour of John Vanderkamp*, University of Toronto Press, Toronto, Canada, 201–222.

The literature reference in the bibliography.

# The journey to metadata 3/6

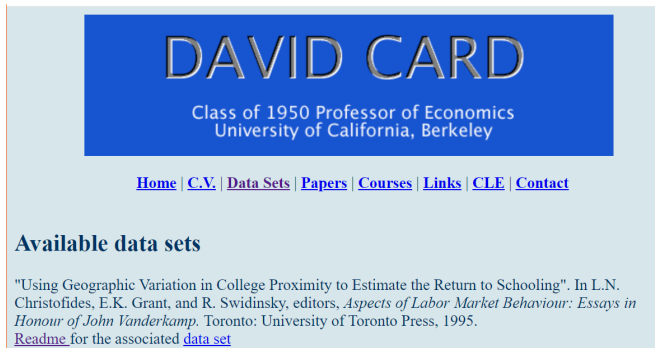


The screenshot shows a Google search interface. The search bar contains the text "david card data". Below the search bar are navigation buttons for "Images", "News", "Videos", "Maps", "Books", "Flights", and "Finance". The search results show "About 262 000 000 results (0,24 seconds)". The first result is from the University of California, Berkeley, with the URL "https://davidcard.berkeley.edu/data\_sets". Below the URL is a link titled "Available data sets" and a snippet of text: "17 Nov 2014 — "Minimum Wages and Employment: A Case Study of the Fast Food Industry in New Jersey and Pennsylvania." (with Alan Krueger), American Economic ...".

Quick and dirty: Google search by author's name.



# The journey to metadata 4/6



The image is a screenshot of a webpage with a light blue background. At the top, there is a dark blue rectangular box containing the name "DAVID CARD" in large, gold, serif capital letters. Below the name, in white text, it reads "Class of 1950 Professor of Economics" and "University of California, Berkeley". Underneath this box, a horizontal line of navigation links is displayed in blue text: "Home | C.V. | Data Sets | Papers | Courses | Links | CLE | Contact". Below the navigation links, the section "Available data sets" is written in bold black text. Under this section, a paragraph of text describes a paper: "Using Geographic Variation in College Proximity to Estimate the Return to Schooling". It lists the authors as L.N. Christofides, E.K. Grant, and R. Swidinsky, and mentions the book "Aspects of Labor Market Behaviour: Essays in Honour of John Vanderkamp" published by the University of Toronto Press in 1995. At the end of the paragraph, there is a blue link labeled "Readme" followed by the text "for the associated data set".

**DAVID CARD**

Class of 1950 Professor of Economics  
University of California, Berkeley

[Home](#) | [C.V.](#) | [Data Sets](#) | [Papers](#) | [Courses](#) | [Links](#) | [CLE](#) | [Contact](#)

**Available data sets**

"Using Geographic Variation in College Proximity to Estimate the Return to Schooling". In L.N. Christofides, E.K. Grant, and R. Swidinsky, editors, *Aspects of Labor Market Behaviour: Essays in Honour of John Vanderkamp*. Toronto: University of Toronto Press, 1995.  
[Readme](#) for the associated [data set](#)

The dataset under the author's homepage.

# The journey to metadata 5/6

```
code_book
Code Book for Ascoli File nls.dat
2
3
4 Note: For more information, see the original article:
5 David Card,
6 "Using Geographic Variation in College Proximity to Estimate the Return to Schooling"
7 NBER Working Paper 4032, August 1994
8
9 This article is published (with the same title and identical tables)
10 in Aspects of Labour Market Behaviour: Essays in Honour of John Vanekamp*
11 edited by Isaac S. Christofides, E. Kenneth Grant, and Robert Swidinsky
12 Toronto: University of Toronto Press, 1995.
13
14 The file contains 3613 observations on men in 1976 cross-section
15 of nls young men (original nls cohort)
16
17 Missing value code is .
18
19 Column Locations      Variable Name and Label
20 -----
21 1 - 3   id      /*sequential id num from 1 to 3229 */
22 7 - 7   year2   /*grew up near 2-yr college*/
23 10 - 10  year3   /*grew up near 4-yr college*/
24 12 - 13  year4   /*grew up near 4-yr public college*/
25 15 - 16  year6   /*grew up near 4-yr priv college*/
26 18 - 19  ed76   /*educ in 1976*/
27 21 - 22  ed66   /*educ in 1966*/
28 24 - 25  age76   /* age in 1976
29 27 - 31  ededu   /*dads education missing=avg*/
30 33 - 33  momedu   /* 1 if dad ed imputed*/
31 35 - 39  momedu   /*mom's education*/
32 41 - 41  momedu   /* 1 if mom ed imputed*/
33 43 - 54  weight  /* nls weight for 1976 cross-section */
34 56 - 56  momeda14 /* 1 if lived with mom and dad age 14+ */
35 58 - 58  sinom14 /*lived with single mom age 14+*/
36 60 - 60  step14 /*lived step parent age 14+ */
37 62 - 62  reg66   /* dummy for region* in 1966 */
38 64 - 64  reg66   /* dummy for region* in 1966 */
39 66 - 66  reg63   /* dummy for region* in 1966 */
40 68 - 68  reg66   /* dummy for region* in 1966 */
41 70 - 70  reg66   /* dummy for region* in 1966 */
42 72 - 72  reg66   /* dummy for region* in 1966 */
43 74 - 74  reg66   /* dummy for region* in 1966 */
44 76 - 76  reg66   /* dummy for region* in 1966 */
45 78 - 78  reg66   /* dummy for region* in 1966 */
46 80 - 80  south66 /*lived in south in 1966*/
47 82 - 82  work76 /* worked in 1976*/
48 84 - 84  work76 /* worked in 1976*/
49 86 - 97  wage76 /*log wage (outliers trimmed) 1976 */
50 99 - 110 wage76 /*log wage in 1976 outliers trimmed */
51 112 - 112 famed /*normalized education class 1-3*/
52 114 - 114 black /* 1 if black*/
53 116 - 116 smm67 /*in smm in 1976*/
54 118 - 118 smm67 /*in smm in 1976*/
55 120 - 120 reg76 /*in south in 1976*/
56 122 - 122 reg78 /*in south in 1978*/
57 124 - 124 reg80 /*in south in 1980*/
58 126 - 126 smm66 /*in smm in 1966*/
59 128 - 132 wage76 /*raw wage cents per hour 1976*/
60 134 - 133 wage78
61 140 - 144 wage80
62 146 - 146 nonint76 /* 1 if noninterview in 76*/
63 150 - 150 nonint80
64 152 - 152 enroll76 /* 1 if enrolled in 76*/
65 154 - 154 enroll80
66 156 - 157 hnw /*raw hnw score*/
67 159 - 161 iq /* a normed iq score*/
68 163 - 163 marr67 /*mar status in 1976 1=marrried, sp. present *
69 165 - 165 marr67
70 167 - 167 marr80
71 169 - 169 libord14 /* 1 if lib oard in home age 14+*/
```

Variable descriptions in the dataset.

# The journey to metadata 6/6

77 List of means, min/max Note some vars are missing for some observations.  
 78 Missing value code is .  
 79  
 80  
 81  
 82

83	Variable	N	Mean	Std Dev	Minimum	Maximum
84	ID	3613	2609.78	1498.51	2.0000000	5225.00
85	NEARC2	3613	0.4317741	0.4953919	0	1.0000000
86	NEARC4	3613	0.6781068	0.4672469	0	1.0000000
87	NEARC4A	3613	0.4921118	0.5000070	0	1.0000000
88	NEARC4B	3613	0.15199950	0.3581845	0	1.0000000
89	ED76	3613	13.2252975	2.7487411	0	18.0000000
90	ED66	3613	10.7428730	2.4590854	0	18.0000000
91	AGE76	3613	28.1752007	3.1718104	24.0000000	34.0000000
92	DAED	3613	10.0028785	3.2960212	0	18.0000000
93	NOADDED	3613	0.2241904	0.4171058	0	1.0000000
94	WOMED	3613	10.3421872	3.2293785	0	18.0000000
95	NOMOMED	3613	0.1143094	0.3182308	0	1.0000000
96	WEIGHT	3613	320318.35	148506.76	75607.00	1782340.00
97	MONDAD14	3613	0.7921395	0.4058326	0	1.0000000
98	SINMOM14	3613	0.1001937	0.3002997	0	1.0000000
99	STEP14	3613	0.0384722	0.1923599	0	1.0000000
100	REG641	3613	0.0448613	0.2063671	0	1.0000000
101	REG642	3613	0.1549958	0.3619508	0	1.0000000
102	REG643	3613	0.1940216	0.3955003	0	1.0000000
103	REG644	3613	0.0691946	0.2588199	0	1.0000000
104	REG645	3613	0.2095212	0.4070232	0	1.0000000
105	REG646	3613	0.0929975	0.2904691	0	1.0000000
106	REG647	3613	0.1101378	0.3131296	0	1.0000000
107	REG648	3613	0.0309982	0.1733394	0	1.0000000
108	REG649	3613	0.0935511	0.2912434	0	1.0000000
109	SOUTH66	3613	0.4126764	0.4923837	0	1.0000000
110	WORK76	3613	0.8350401	0.3711987	0	1.0000000
111	WORK78	3613	0.7951232	0.4413287	0	1.0000000
112	LMAGE76	3010	6.2615319	0.4437977	4.4681762	7.7848993
113	LMAGE78	2459	6.3291080	0.4442450	4.6965200	8.2409240
114	FAMED	3613	5.9128148	2.6504318	1.0000000	9.0000000
115	BLACK	3613	0.2300028	0.4208925	0	1.0000000
116	SMSA76R	3613	0.6947135	0.4605924	0	1.0000000
117	SMSA78R	3319	0.6929798	0.4613273	0	1.0000000
118	REG76R	3613	0.3996079	0.4898978	0	1.0000000
119	REG78R	3319	0.3965043	0.4893059	0	1.0000000
120	REG90R	3227	0.4028509	0.4905473	0	1.0000000
121	SMSA66R	3613	0.6426792	0.4792768	0	1.0000000
122	WAGE76	3017	576.088300	263.8199090	25.0000000	2404.00
123	WAGE78	2656	724.5591114	526.1991520	17.0000000	17628.00
124	WAGE90	2520	869.8940476	492.1729068	27.0000000	13857.00
125	WOMINT8	3613	0.0319718	0.1724447	0	1.0000000
126	WOMINT80	3613	0.1306364	0.3089479	0	1.0000000
127	ENROLL76	3613	0.0946582	0.2927827	0	1.0000000
128	ENROLL78	3317	0.0654206	0.2473038	0	1.0000000
129	ENROLL80	3220	0.0589851	0.2345066	0	1.0000000
130	HMW	3543	33.4891335	8.6918079	0	56.0000000
131	TQ	2470	102.5578543	15.4455073	50.0000000	156.0000000
132	MARSTA76	3604	2.3871032	2.1096377	1.0000000	6.0000000
133	MARSTA78	3319	2.2136186	2.0058342	1.0000000	6.0000000
134	MARSTA90	3227	2.1041215	1.9088835	1.0000000	6.0000000
135	LIBCRD14	3598	0.6717621	0.4686972	0	1.0000000
136						

Variable means in the dataset.

# Missing values

- ▶ Missing values in data do not inherently mean the data is corrupted.
- ▶ Example: what is the annual wage for a person who is not employed? Is zero a meaningful value for their wage?
- ▶ By default, R's `lm` ignores any rows where any regressors (or the dependent variable) have missing values.

# Missing values

- ▶ There are at least three ways to cope with missing data. The appropriate method depends on the context.
  - ▶ Dropping the rows with missing data (you can use this during this course; it's what R does by default)
  - ▶ Multiple imputation (replace missing values with predictions); not covered in this course
  - ▶ Discretizing continuous values and including missingness as an additional category; not covered in this course

## Dropping missing values in R

```
schooling <- read.table(  
  'schooling.txt', header = T)  
schooling <- subset(schooling, !is.na(iqscore))
```

### 3. Instrumental variables

- ▶ Instrumental variables are one way of dealing with cases where the regular OLS assumptions don't hold
- ▶ It is one (of many) causal identification strategies; others include
- ▶ Other popular strategies:
  - ▶ Differences-in-differences (DiD or DD)
  - ▶ Regression discontinuity (RD)
  - ▶ Randomized controlled trials (RCT)
  - ▶ Adjustment/unconfoundedness (see also: matching and weighting)

# Further stuff on causality

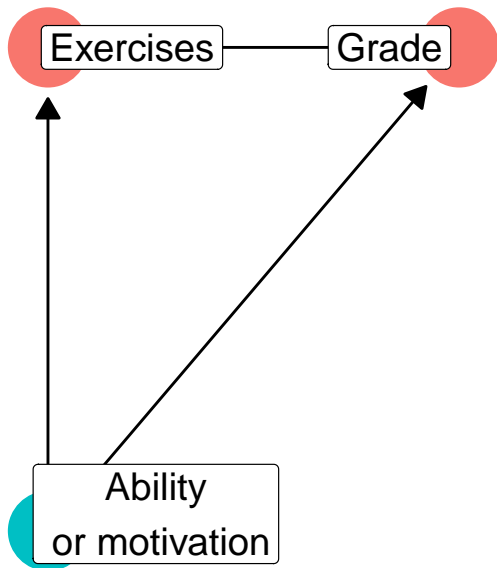
- ▶ The books: Angrist and Pischke: Mostly Harmless Econometrics, or Mastering Metrics (less technical)
- ▶ Complementary reading
  - ▶ Pearl: The Book of Why
  - ▶ Imbens 2020: Potential Outcome and Directed Acyclic Graph  
Approaches to Causality: Relevance for Empirical Practice in Economics
- ▶ Courses
  - ▶ Applied Microeconometrics I and II (Aalto University) cover empirical research using micro-level data
  - ▶ Applied Macroeconometrics I and II (UoH) cover time series using macroeconomic data
  - ▶ Applied courses are very helpful if you want to do or at least cite and understand causal empirical work in your thesis (you probably will)

# The problem to be solved

- ▶ Suppose we want to know how forcing everyone to do exercises affects learning
  - ▶ We also assume grades are a good measure of learning
- ▶ But those scoring well on exercises probably were more motivated or able to start with
  - ▶ They might have done well on the exam without the exercises
  - ▶ Motivation/ability is usually (partially) unobserved
- ▶ Thus, the regression might show positive correlation between exercises and course grade even if exercises do nothing for learning (more commonly, the regression might under- overestimate the effect)



# The problem to be solved



colour



Observed

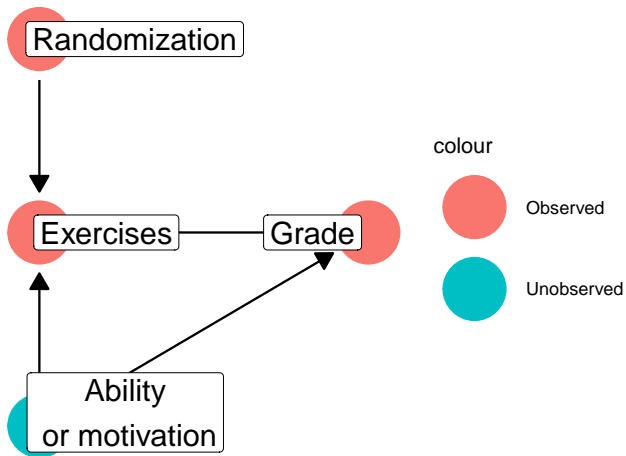


Unobserved

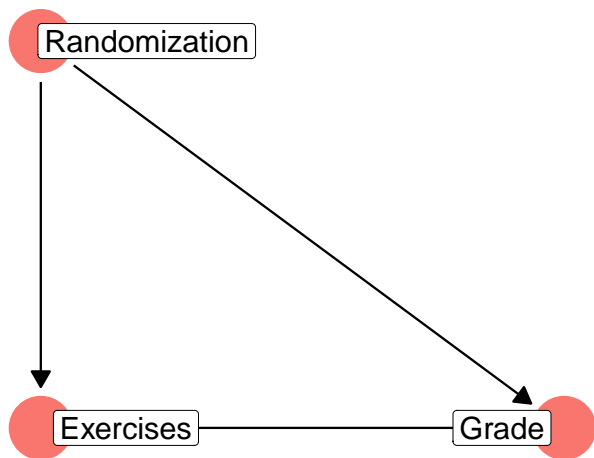
# The solution

- ▶ One solution is a randomized controlled trial (RCT)
- ▶ A randomly selected half has to do the exercises to pass the course
- ▶ Successful random assignment is, by construction, independent of grades, ability, and whatever else you might think of.
- ▶ However, the RCT might have imperfect takeup: not everyone assigned to do exercises does them
- ▶ In this case, we can (under certain conditions) still use random assignment as an instrument for the exercises

# The solution



## Bad instrument: endogenous



The instrument should not be correlated with the outcomes.

# Bad instrument: irrelevant

Randomization

A diagram illustrating an instrument. A red circle is partially behind a white box containing the text 'Randomization'. This box is positioned above the 'Exercises' box in the diagram below.

Exercises

A diagram showing a causal relationship. A horizontal line connects a white box labeled 'Exercises' on the left to a white box labeled 'Grade' on the right. Both boxes are partially overlaid by red circles.

Grade

The instrument should be correlated with the regressor.

# Instruments: only one strategy

- ▶ Instrument relevance is easy to test for, but a difficult problem to solve when it occurs (and also sadly common)
- ▶ An even bigger problem is that the endogeneity assumption is only partly testable
- ▶ The Sargan test **assumes that at least some instruments are valid.**
  - ▶ See e.g. the textbook p. 168
- ▶ IV results are local: they are informative about how instrument-induced changes in  $X$  affect  $Y$
- ▶ Any identification strategy (other than the idealized RCT) will include some untestable assumptions; generally, the question is about the plausibility of the assumptions

# Instruments: only one strategy

- ▶ Applied work based on IV alone is arguably more prone to criticism than, for example, DiD and RD (which have their own issues)
- ▶ Coming up with such criticism is arguably the easiest job in economics
- ▶ The extreme (silly) version of a skeptical take is known as "Friends don't let Friends do IV"
- ▶ A more sane approach comes under the heading "Friends \*Do\* Let Friends Do IV":

*If you are going to use an observational IV, you do need to think very carefully. — Identifying causal effects is hard. Willingly limiting yourself to a subset of methods and declaring one method off-limits is like a football coach saying he doesn't want his quarterback to ever try to pass the ball.*