Paper --see guidelines

- attempt at a coherent regression analysis.

- start with a question or hypothesis
  sharper the better
  doesn't have to be supported, just
  well tested.

- any data -- the class data are fine
  I will work with you on this.

- structure - walk through handout.

Diff-in-diff

Using regression for causal analysis    RD, diff-in-diff,
                                          fixed effects

  Today diff-in-diff to take advantage of natural experiments
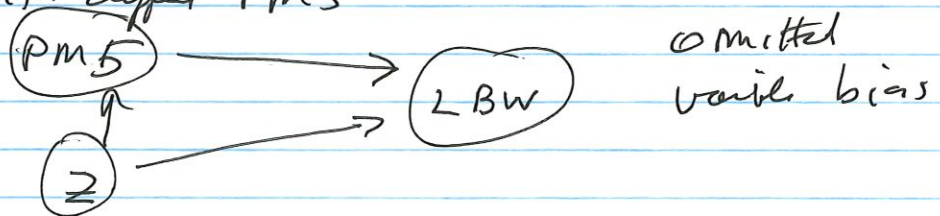
Recall the basic idea of natural experiments

  NOT:    Data set ⟶   what can I do
                        with this data set?

  BUT:  Event that
        creates exogenous
        variation in a        ⟶    Assemble data
        measure/policy of            and construct model
        interest                     to do the test
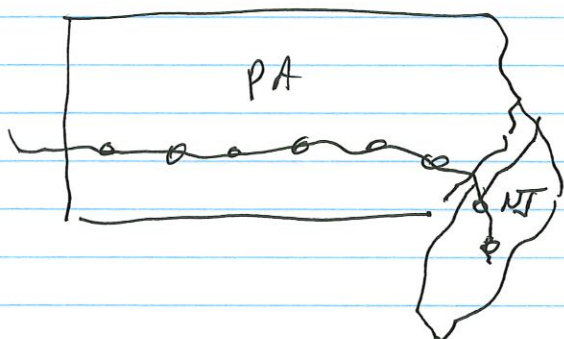
  Example:  Currie and Walker:  what is the effect
            of ~~air~~ auto exhaust pollution on child
            health

  NOT:    Assemble data on, say, PM5 and
          ~~the~~ LBW and correlate
          LBW = a + b₁ PM5 + controls

          because: why do people live in different areas
          with different PM5



PM5 ⟶ LBW          omitted
                    variable bias

She uses the natural experiment of the roll-out
of EZPass on the NJ & PA turnpikes



lots of TOLL booths.  originally everyone stopped, then
        EZPass -- no slowing down. ⟹ Huge reduction
    in emissions

If emissions matter, then you should observe reduction
    in the bad child outcomes for people living closest
    the the toll booths.

So you want to compare birth outcomes
    ~~w/birth~~
        1. Before and after the change to EZPass
        2. For women living very close us. a bit
            bit further way from them
                        < 2 km  us  2-5 km

take only births < 5 km of a toll booth

and run $BW = a + b_1 \, EZ \, pass + b_2 < 2km$

$\qquad\qquad + b_3 \, EZ \, pass + b_2 \, km$

so: $b_1$ is the main effect of EZ pass

$\quad b_2$ is the main effect of < 2km relative to 2-5

$\quad b_3$ is the extra effect of EZ pass for the < 2km ers.

$\qquad\qquad$ (relative to the 2-5 ers)

Go over table 3    $\hookrightarrow BW \, coeff = -.0093$

$\qquad$ base rate (in Table 1) is abt 10% (.10)

$\qquad$ so reduction is abt 10% / pat $= \sim 10\%$

Actual model is more sophisticated because it
includes multiple births to the same mother.

so add another difference

$\quad$ 3. Among births to the same women.

$\qquad$ Fixed effects

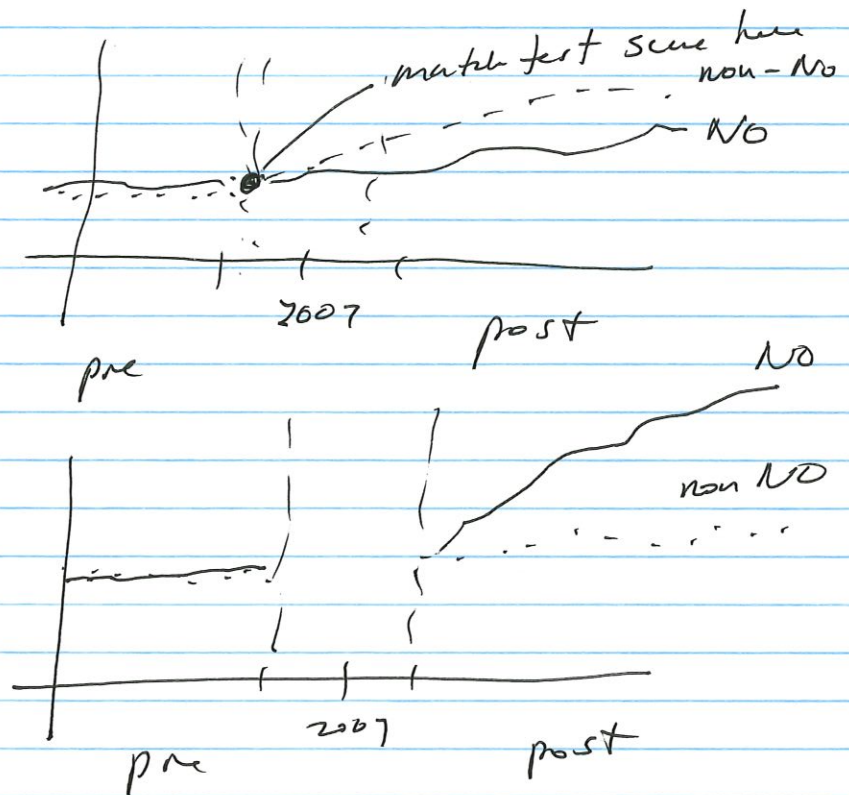One other example. Optional -- impact of charter
schools in New Orleans.

Hurric Katrina hit in 2007 and NO
school district became nearly 100% charter schools.
Did kids do better?

Do better than whom?

- All other kids in US? in LA?
No -- than kids in neighboring school districts
with similar pre-Katrina test scores.

Assemble data on

Test score $= a + b_1 \, 2004 + b_2 \, 2005 + \cdots b_n \, \frac{2014}{~~~~}$

$+ c_1 \cancel{*} \text{white NO} + d_1 \, NO * 2004 + d_2 \, NO \cdot 2005$

$\cdots + d_n \, 2014$

Expect that $d_1, d_2$ would be close to zero

but d's for 2009-2014 ~~would~~ would be
either positive or negative.

Applied Regression Paper Format Suggestions

In general, I am hoping that you will be spending almost as much time writing the paper as analyzing the data in it. Too often students spend almost all of their time futzing with their analyses and very little time writing it up. Our PhD program wants to develop both your analytic and writing skills.

1.) Abstract: When I sit down to write a paper, I force myself to write a ~150-word abstract first, because it you can't write the abstract, you don't yet have a story to tell. Every paper tells a story, with motivation for why what you are examining matter, your research questions, your methods and your results. You should be able to say that in 150 words.

2.) Introduction: Here you want to set up your hypothesis. In a few paragraphs, explain why your research is interesting and needed, your research question/hypothesis, and very briefly discuss any background research or theories that inform your research question (and I not expecting as much here as I would in a second-year paper).

3.) Describe your methods: sample, measures, and study procedures. Your explanation should be clear and complete enough so that I could replicate what you did if I had access to your data. Use tables of descriptive statistics (means, standard deviations, proportions, with your key measures at the top of the table) within this section. If you are looking at differences between groups, do descriptive statistics by those groups and provide p-levels from statistical tests of group differences. Describe all your independent variables in this way.
Codes to use:
tab, sum, estout
*you may need to create dummy variables for your groups. Use codebook in stata and/or the data documentation to find the appropriate identifiers for your dummy variables (Example in HW 2). Take a look at some of our course readings for ways of doing this concisely.

4.) Describe the kind of analyses you will use to test your hypotheses. Frame your description in light of the types of analyses we have discussed in class: OLS, logistic, spline, quadratic, fixed effects etc. (also residualized change and simple change, which we will cover in future classes). Including formulas may be useful.
How will you handle missing data? Are you limiting your sample in any way? How does the limited sample compare to the entire sample (use descriptives). Are your variables highly correlated? Will this correlation present a problem (multicollinearity).

5.) Present your key results in text and tables, identifying which of your coefficients tests your key hypotheses. Your tables should be publication-ready: in APA or another format with variable descriptors that make sense (not the actual names of your variables within stata). In text, report the associations between variables as we have in class and lab—what do the numbers mean?

6.) Briefly summarize your results in light of your question, hypothesis and the theories or background that inform your paper. What are the implications for your results? Are there any limitations to your study (there should be)? I am expecting this to be much shorter than would usually be the case in an article-type paper.

Continuing with our look at using regression to estimate causal relationships

Fixed effects    --     adding dummy variables for all "units" in the sample

       e.g. all schools in the sample

         " families " " sample

         " sites " a multi-site experiment.

"Econometric" fixed effects    as opposed to HLM "fixed effects" which are fixed (vs. random effects"

      Two broad uses

1. Show where the action is. e.g. Variance decomposition

2. Powerful control for omitted variable bias
     (Powerful because it controls for both measurable AND unmeasurable sources of bias)

Heresy 101: (Econometric) fixed effects adjustments are a powerful and under appreciated technique for reducing bias

      (rare in development studies)

Show where the section is

Fryer and Levitt    Table 2

Fall K Math

Black        −.663 (.025)

Hispanic     −.738 (.024)

⋮

Constant     .307 (.013)

$R^2$      .11  ⟵  11% of variation in test score
                  is accounted for with 4 dummy
                  variables

Highly significant in a statistical sense __but__

~~there~~ ~ 90% of the variation is __within__

rather than across race/ethnic groups.

But now look at Table 7: Does school quality
                    Explain Black Students'
                              losing ground?

change in gap K to 3rd

| | | |
|---|---|---|
| Math | −.243 | −.180 |
| | (.052) | (.061) |
| Reading | −.343 | −.214 |
| | (.057) | (.065) |
| include school fixed effects? | No | Yes |

-.243 vs. -.180  => 3/4 of the growth in the gap occurs within schools rather than between schools

Most of the action is not caused by the fact that Black and White students attend different schools
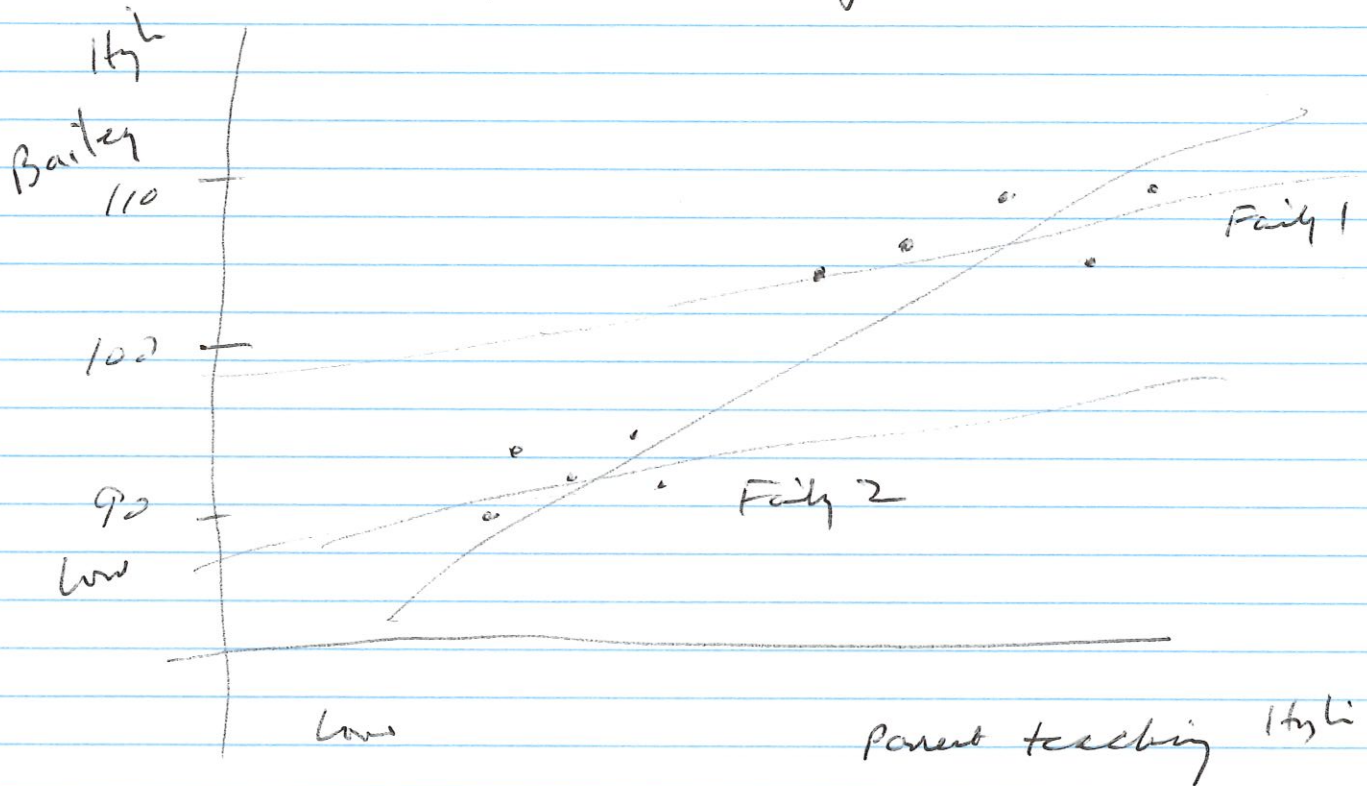
Not always the case
      Michelle and lead in water
Across Portland substantially higher lead in school water for Black students
within or across the schools
with or across class rooms within schools

Fixed effects as a way of reducing omitted variable bias.

Suppose you are investigating the relationship between parent cognitive stimulation and child cognitive development

NCATS - print teaching scale at 9 months
interms    Bailey Mental Scale   at 2 years
           Behavior Rating Scale  at 2 years.



Suppose prints are elevated at 9 months with all of their child

Biased across but not within

$$\text{Baily} = a + b_1 \text{ NCATS} + \underset{\text{Faly 1}}{\text{whh}} + \ldots$$

allows for sepret. intercept <del>within taxa</del> for each unit

Note:   no whether & NCATS intrectn
Faly 1

slopes are assumed to be the same
within all taxa.

if you have 1000 taxa you don't want
1,000 interactn terms
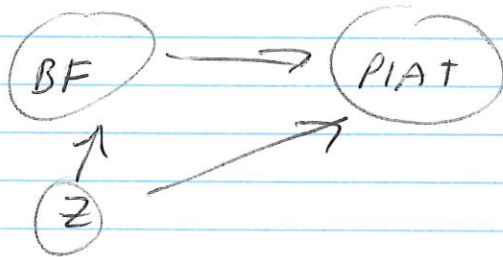
Cassie's results for AlCATB
vs results for Birthright

Fixed effects models

(1) Add dummy variable for every unit

(2) Use xtreg and specify unit

(3) Transform all variables so that they are deviations from their unit means.

Why does that work.

$$PIAT \ math = a + b_1 \ Breast \ Fed + b_2 \ \overset{Fixed}{Family \ chars}$$

$$+ b_3 \ other \ stuff$$



Suppose sibling data, 2 sibs A:B

$$PIAT_A = a + b_1 BF_A + b_2 Fam_A + b_3 OS_A$$

$$PIAT_B = a + b_1 BF_B + b_2 Fam_B + b_3 OS_B$$

but $Fam_A$ is the same as $Fam_B$

$$PIAT_B - PIAT_A = a + b_1 (BF_B - BF_A) + O$$

$$+ b_3 (OS_B - OS_A)$$

$b_1$'s ~~are the~~ have the same integrity

Go to  Table 4

Big reduction. Should earn au layer
                             a cost 1

Go back to | Curie and Walker |