

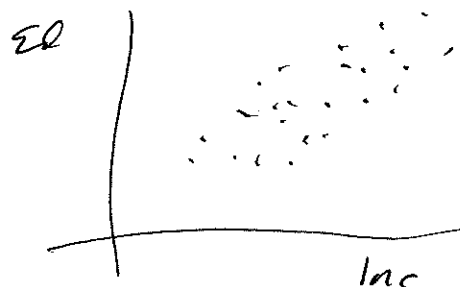
5th class Part I

- 1 -

Logistic regression -- very annoying but necessary

Thus far, our dependent variables have been continuous

What if dichotomous?
(two valued)



$$\rightarrow y = a + b_1 x$$

First point: variation in a dependent variable is a terrible thing to waste!

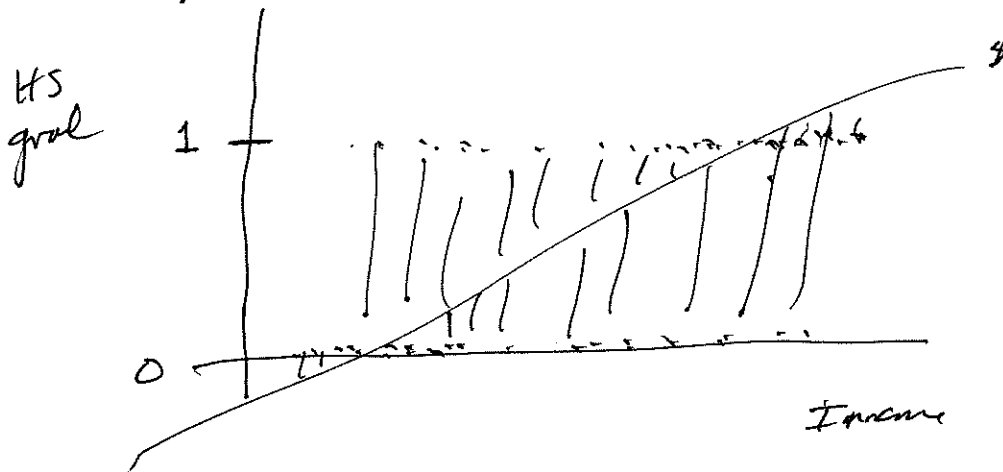
Deaton et al. analyzed year of Ed
but also whether he grad^d (1,0)

In most cases, go with the continuous variable even
if it is not well distributed

If points of the distribution are of particular interest, (1,0)
or quantile regression

But sometimes you are stuck with a 1,0 variable.

Suppose $(1,0) = a + b, \text{ income}$ ^{Fairly}
 HS grad



What does ~~scatter~~ scatter look like

OLS: do your job! minimize squared residuals

It gives you the best fitting slope -- say .03
 (.01)

How do interpret .03

Second big point: think of the 0,1 scale as
 a 0 → 100 percent chance of being a 1

If ^{exactly} $\frac{1}{2}$ of people, like me finished high school,
 then my chance of finishing are 50%

If $\frac{1}{4} \rightarrow 25\%$

$\frac{3}{4} \rightarrow 75\%$

Now the .03 makes a lot of sense
(.01)

-- it is the increase in probability of graduating
high school associated with a \$10,000 increase
in income

So if I am at 50% at the .03 is a causal
estimate, then my chances will be 53% with
more money

Very straightforward! $1.0 = a + b$, Func
is a linear probability model

⇒ changes in ~~prob~~ probability are a
linear function of income

.25	→	.28
.75	→	.78

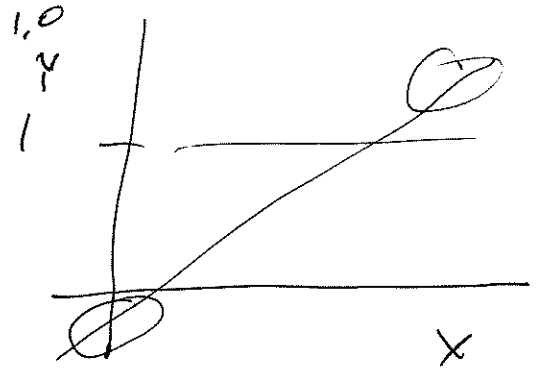
Lot of people use LPMs because ~~you can do~~
they aren't bad until \bar{Y} is between .2 and .8

But LPMs offend the sensibilities of statisticians!

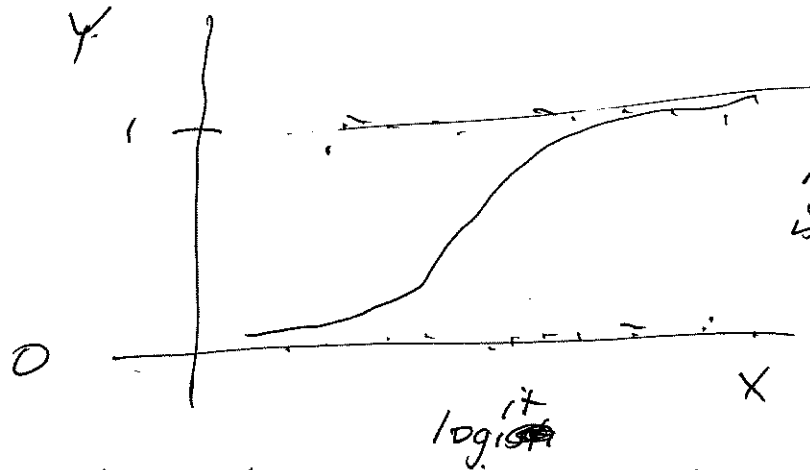
Horror!

LPN problems:

- ① Predictive behavior above 1
- ② Non constant variance



Want a nonlinear model that doesn't exceed bounds



Note that slope is constantly changing

two ~~word~~ models -- ~~logistic~~ and probit

Almost interchangeable, Logit is more popular

Instead of $P = a + b, \text{ For linear}$

$$\ln\left(\frac{P}{1-P}\right) = a + b, \text{ For linear}$$

ugh!

5th class Part I

-5-

$\ln\left(\frac{p}{1-p}\right)$
 ↑ natural logarithm
 ↑ odds ratio

Odds $\frac{p}{1-p}$

$\frac{p}{1-p}$ odds
 $\frac{.5}{.5} = 1$

$\frac{.67}{.33} = 2$
 even odds
 $2 + 1$

If we run $\ln\frac{p}{1-p} = a + b_1$, then $\ln e$

$\frac{.33}{.67} = \frac{1}{2}$
 $1 \text{ to } 2$

what is value of b_1

literally change in log odds associated with a one unit \$10K change in fixed income

Still unclear -- how to talk about change in log odds??

$$\ln\left(\frac{p}{1-p}\right) = a + b_1 \text{, Fixed Inc}$$

because $e^{\ln x} = x$

$$e^{\ln \frac{p}{1-p}} = \frac{p}{1-p} = e^{a + b_1 \text{, Fixed Inc}}$$

if x change by 1 unit $\Rightarrow \Delta \frac{p}{1-p} = e^{b_1}$

so take logit coeff and exponentiate to get change in log odds

5th class Part I

In Duacac et al. .23
(07)

-6-
 $.00 \Rightarrow e^0 = 1$
1 = no effect on odds

literally Δ in log odds of complexity is associated
with an extra \$10K in income

If log odds change by .23 the odds ratio
changes by $e^{.23} = 1.259$

so odds change by 25.9%

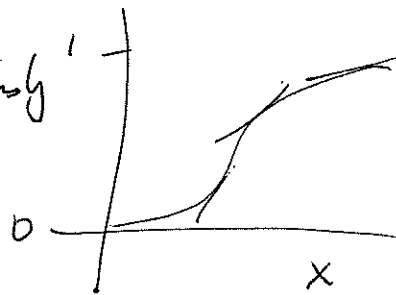
if -.15 then $e^{-.15} = .86$

odds drop by 14%

most medical articles focus on odds ratios

Get me back to changes in probabilities!!

Key challenge: Δ prob is changing continuously



Need some way of evaluating slope
at some relevant point on the curve

"Margins" give you a probability intercept at
mean of X's or any other place you
want.

5th class Part I

7
-6-

My preference is to have regression tables with coefficients and standard errors that have probably interpretation

Key takeaways

1. ~~the~~ LPM's are usually fine for exploration and econ journals
2. logit with "margins" is good because of its probability interpretation
3. log-odds and odds ratios if you must!

Compare Prospecing = $a + b_1 TV$

	coefficient	s.e.	t or z ratio
OLS (LPM)	.066	.015	4.52
log odds	.405	.142	4.27
odds ratio ($e^{.405}$)	1.500	.094	4.27
margins at mean	.089	.021	4.27

$\beta(P)(1-P)$
↑
margin mean

This week marks a turning point in the class

Up to Now: basic tools of regression analysis

Next 2⁺ weeks: using regression for causal analysis

"Every regression is an experiment". Q: is it a convincing causal experiment?

How to identify causal effects?

also "identification strategies"

use example of effect of music on ~~parent~~ child ed. ^{5K more}

<u>Procedure</u>	<u>Strengths of identification strategy</u>	<u>Problems</u>
1. Random assignment (of 5K)	Very strong data	Imbalance, generalizability
2. Regression with controls	Often weak	Omitted variable bias

3. Quasi-experimental approaches "natural experiments"

- ~~one~~ regression discontinuity -- Wang + Cortes

- difference in difference -- Akee, Dahl & Lochner

based on some exogenous change

in variable of interest

→ sibling models

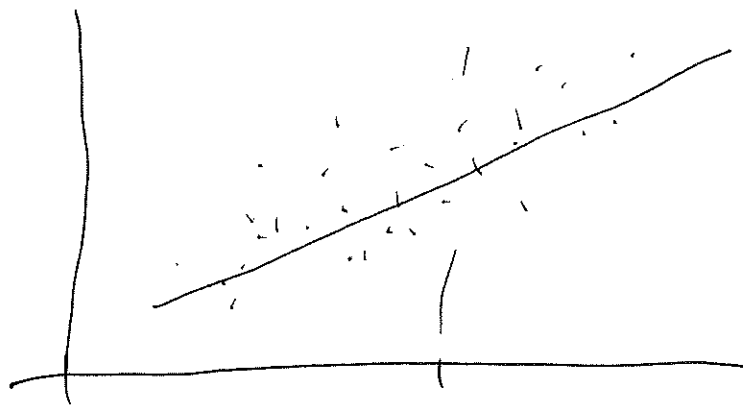
Week 6 is about regression discontinuity

- Did some interesting event description what would normally be an unremovable correlation?

e.g. Double dose algebra for ninth-grade student

In the absence of double dose --

9th grade math

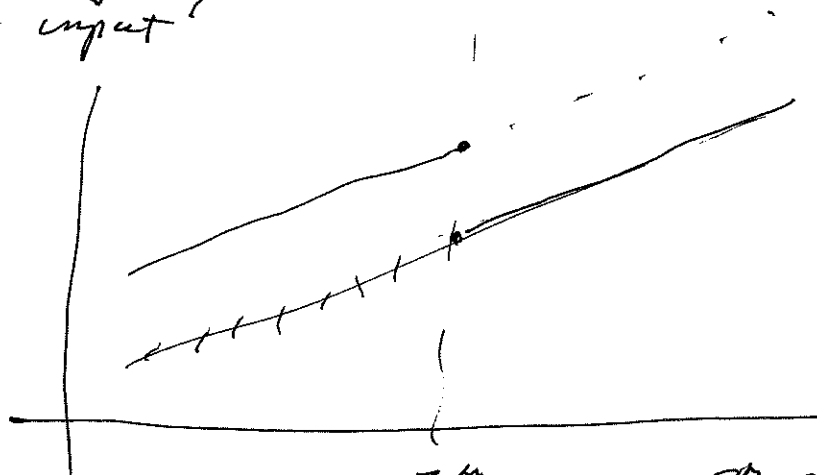


50th percentile

8th grade math

DD targets students below the 50th percentile in 8th grade math input?

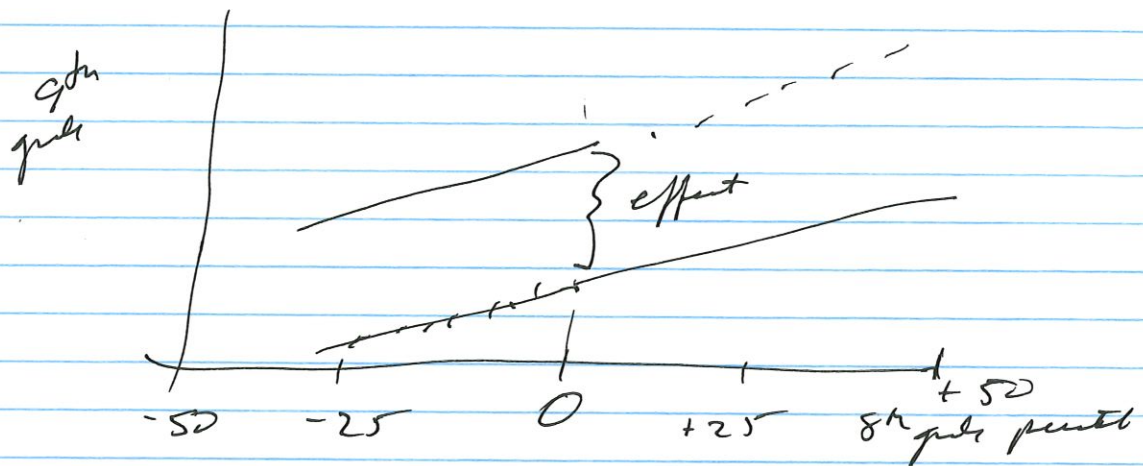
9th grade math



50th

8th grade math

How to model DD effects. First, center data on the 50th percentile



$$9^{\text{th}} \text{ grade math} = a + b_1 \text{ 8th grade percentile} + b_2 \text{ whether } < 50^{\text{th}} \text{ percentile}$$

For above 50th percentiles:

$$9^{\text{th}} \text{ grade} = a + b_1 \text{ 8th grade percentile}$$

for below 50th percentiles:

$$9^{\text{th}} \text{ grade} = (a + b_2) + b_1 \text{ 8th grade percentile}$$

b_2 is the intercept

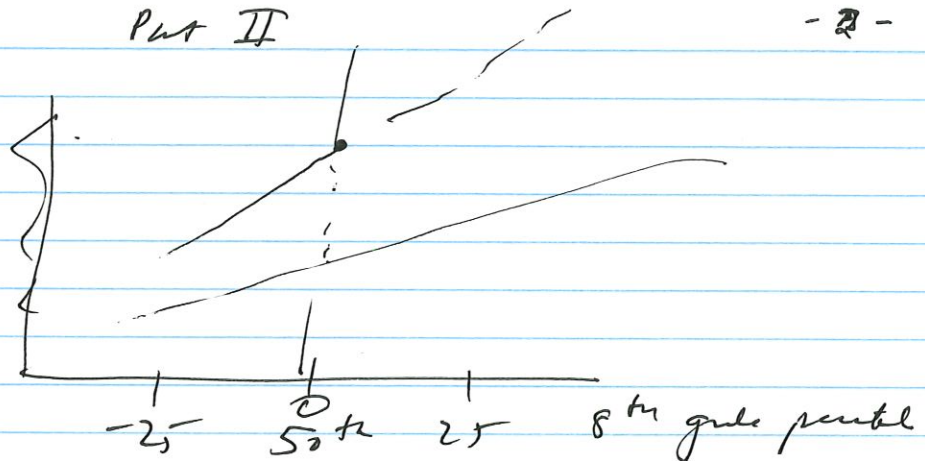
because intercept is evaluated at 50th percentile

5th class

Part II

4
-2-

maybe
slopes
are
different



$$9^{\text{th}} \text{ grade} = a + b_1 8^{\text{th}} \text{ grade} + b_2 \text{ water} + b_3 \text{ water}$$

< 50th < 50th
8th grade 8th grade

for > 50th percentile

$$9^{\text{th}} = a + b_1 8^{\text{th}}$$

for < 50 percentile

$$9^{\text{th}} = (a + b_2) + (b_1 + b_3) 8^{\text{th}} \text{ grade}$$

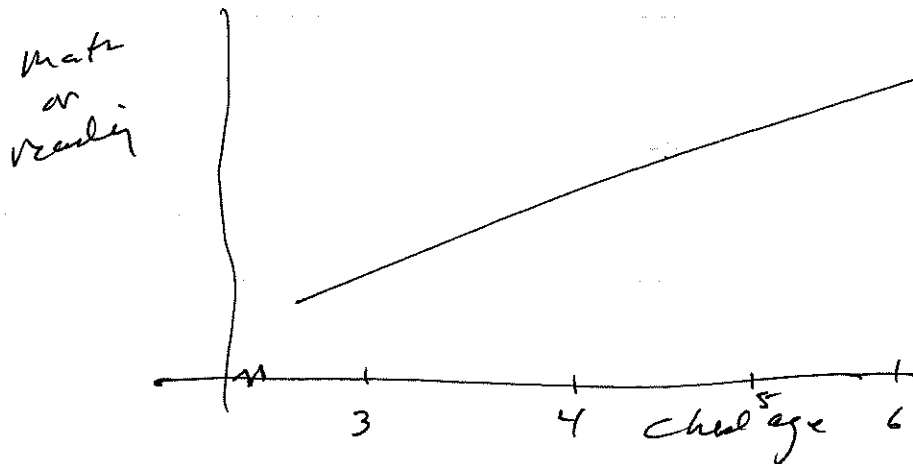
Allowing for different slopes enables you to
better nail down the jump at the 50th
percentile.

Week 6 is about regression discontinuity -

Did some interesting event disrupt what would normally be a ~~more~~ uncorrelated correlation?

offer

e.g. Preschool disrupts the relationship between child age and, say, math scores



Now offer preschool to ~~the~~ on 4th birthday and results

