

Lecture 8 - Beyond basics of OLS.pdf

Thursday, October 8, 2020 3:59 PM



Lecture 8 -
Beyond...

Beyond the basic of OLS

Mauricio Romero

1

Beyond the basic of OLS

A few things that don't get enough attention

Error structure

Statistical power

2

Beyond the basic of OLS

A few things that don't get enough attention

Error structure

Statistical power

3

An example

- A regression of wages on: Age (in years), race (black=1) and IQ percentile (0-100)
- For every year, we expect wages to change by $\widehat{\beta}_{age}$ USD
- On average, we expect wages to higher/lower for blacks by $\widehat{\beta}_{black}$ USD than for non-blacks
- For every **percentage point** increase in IQ, we expect wages to change by $\widehat{\beta}_{IQ}$ USD

8

Simulations!

```
library(wooldridge)
library(stargazer)
data("wage2")
wage2$IQ_Percentile=quantile(wage2$IQ, seq(0, 1, 0.1))
levlev=lm(wage ~ IQ_Percentile + age + black, data = wage2)
summary(levlev)
stargazer(levlev, title="Level-Level", align=TRUE,
  type="latex", omit.table.layout="la",
  out="Lectures/tables/levlev.tex",
  covariate.labels=c("IQ (percentile)", "Age", "Black(=1)"),
  digits=2, digits.extra=1, no.space=T, colnames=F,
  dep.var.caption="", dep.var.labels="Wage",
  column.sep.width="0pt", headers=F,
  omit.stat=c("adj.rsq", "rsq", "f", "ser"))
```

9

```
Call:
lm(formula = wage ~ IQ_Percentile + age + black, data = wage2)

Residuals:
    Min       1Q   Median       3Q      Max
-803.60 -271.87  -62.62  212.27 2174.38

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  332.4888   150.2711    2.213  0.0272 *
IQ_Percentile  0.1320    0.5615    0.235  0.8141
age          19.4666    4.1241    4.720 2.72e-06 ***
black       -248.0806   38.2995  -6.477 1.51e-10 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 391.2 on 931 degrees of freedom
Multiple R-squared:  0.06681, Adjusted R-squared:  0.0638
F-statistic: 22.22 on 3 and 931 DF, p-value: 6.71e-14
```

10

BLACKASS

Level-Level	
	Wage
IQ (percentile)	0.13 (0.56)
Age	19.47*** (4.12)
Black(=1)	-248.08*** (38.30)
Constant	332.49** (150.27)
Observations	935
Note:	*p<0.1; **p<0.05; ***p<0.01

*SALARIO (IQ=0, Year=0)
BLANCO*

$$y = \beta_0 IQ + \beta_1 AGE + \beta_2 BLACK + \epsilon$$

11

Log-level Regression

- If you have a log-level regression

$$\ln(y_i) = \beta_0 + \beta_1 x_i + u_i$$

- If you increase x by one, we expect y to change by $100\beta_1$ percent
 - Technically $\% \Delta y = 100(e^{\beta_1} - 1)$
 - But $\% \Delta y = 100(e^{\beta_1} - 1) \approx 100\beta_1$ for values $-0.1 < \beta_1 < 0.1$
- You can only include observations for which $y_i > 0$
- Only do it if this doesn't introduce bias into your sample
 - In general, only do it if $y_i > 0$ for almost all i
 - Adding 1 or 0.1, or 100 is not a valid fix

$$e^{\beta_1} - 1 \approx \beta_1$$

$$-0.1 \leq \beta_1 < 0.1$$

$$Y_i = \text{SALARIES} \in [0, \infty] \rightarrow$$

$$\text{Log}(Y_i + 0.001)$$

$$\text{Log}(Y_i + 0.1)$$

$$\text{Log}(Y_i + 1)$$

$$\text{Log}(Y_i + 10)$$

12

An example

- A regression of $\ln(\text{wages})$ on: Age (in years), race ($\text{black}=1$) and IQ percentile (0-100)
- For every year, we expect wages to change by $100\widehat{\beta}_{\text{age}}$ percent
- On average, we expect wages to be higher/lower for blacks by $100\widehat{\beta}_{\text{black}}$ percent than for non-blacks
- For every percentage point increase in IQ, we expect wages to change by $100\widehat{\beta}_{\text{IQ}}$ percent

13

Simulations!

```
loglev=lm(log(wage) ~ IQ_Percentile + age + black, data = wage2)
summary(loglev)
stargazer(loglev, title="Log-Level", align=TRUE,
           type="latex", omit.table.layout="la",
           out="Lectures/tables/loglev.tex",
           covariate.labels=c("IQ (percentile)", "Age", "Black(=1)"),
           digits=2, digits.extra=1, no.space=T, colnames=F,
           dep.var.caption="", dep.var.labels="ln(Wage)",
           column.sep.width="0pt", headers=F,
           omit.stat=c("adj.rsq", "rsq", "F", "ser"))
```

14

```
Call:
lm(formula = lwage ~ IQ_Percentile + age + black, data = wage2)

Residuals:
    Min       1Q   Median       3Q      Max
-1.98581 -0.25765  0.01094  0.27996  1.30084

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  6.128e+00  1.556e-01  39.378 < 2e-16 ***
IQ_Percentile -1.153e-05  5.814e-04  -0.020  0.984
age           2.083e-02  4.271e-03   4.878 1.26e-06 ***
black        -2.852e-01  3.966e-02  -7.191 1.33e-12 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.4052 on 931 degrees of freedom
Multiple R-squared:  0.07746, Adjusted R-squared:  0.07449
F-statistic: 26.06 on 3 and 931 DF, p-value: 3.438e-16
```

15

Log-Level

	ln(Wage)
IQ (percentile)	-0.000 (0.001)
Age	0.02*** (0.004)
Black(=1)	-0.29*** (0.04)
Constant	6.13*** (0.16)
Observations	935
Note:	*p<0.1; **p<0.05; ***p<0.01

16

Level-log Regression

- If you have a **Level-Log** regression

$$y_i = \beta_0 + \beta_1 \ln(x_i) + u_i$$

- If you increase x by one percent (**NOT BY ONE PERCENTAGE POINT!**), we expect y to change by $\frac{\beta_1}{100}$ units of y
- You can only include observations for which $x_i > 0$
- Only do it if this doesn't introduce bias into your sample
 - In general, only do it if $x_i > 0$ for almost all i
 - Adding 1 or 0.1, or 100 is not a valid fix**

17

An example

- A regression of *wages* on: $\ln(\text{Age})$, race ($\text{black}=1$) and $\ln(\text{IQ})$ (IQ is the percentile)
- For an increase in 1 percent in age, we expect *wages* to change by $\frac{\beta_{\text{Age}}}{100}$ USD
- On average, we expect *wages* to be higher/lower for blacks by β_{black} USD than for non-blacks
- For an increase in 1 percent in the IQ percentile (that is, a percent change in percentage points), we expect *wages* to change by $\frac{\beta_{\text{IQ}}}{100}$ USD

18

Simulations!

```
levlog=lm(wage ~ log(IQ.Percentile) + log(age) + black, data = wage2)
summary(levlog)
stargazer(levlog, title="Level-Log", align=TRUE,
  type="latex", omit.table.layout="la",
  out="Lectures/tables/levlog.tex",
  covariate.labels=c("ln(IQ (percentile))", "ln(Age)", "Black(=1)"),
  digits=2, digits.extra=1, no.space=T, colnames=F,
  dep.var.caption="", dep.var.labels="Wage",
  column.sep.width="0pt", header=F,
  omit.stat=c("adj.rsq", "rsq", "F", "ser"))
```

19

```

Call:
lm(formula = wage ~ log(iq_Percentile) + log(age) + black, data = wage2)

Residuals:
    Min       1Q   Median       3Q      Max
-803.07 -271.50 -60.65  210.88 2180.13

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  -1340.20    536.08  -2.500   0.0126 *
log(iq_Percentile)   13.98     49.60   0.282   0.7782
log(age)         648.41    136.38   4.754 2.30e-06 ***
black           -247.97     38.29  -6.477 1.51e-10 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 391.2 on 931 degrees of freedom
Multiple R-squared:  0.06713,    Adjusted R-squared:  0.06412
F-statistic: 22.33 on 3 and 931 DF,  p-value: 5.731e-14

```

20

Level-Log

	Wage
<u>ln(IQ (percentile))</u>	13.98 (49.60)
<u>ln(Age)</u>	648.41*** (136.38)
<u>Black(=1)</u>	-247.97*** (38.29)
Constant	-1,340.20** (536.08)
Observations	935
Note:	*p<0.1; **p<0.05; ***p<0.01

21

log-log Regression

- If you have a log-level regression

$$\ln(y_i) = \beta_0 + \beta_1 \ln(x_i) + u_i$$

- If you increase x by one percent (**NOT BY ONE PERCENTAGE POINT!**), we expect y to change by β_1 percent
- You can only include observations for which $x_i > 0$ and $y_i > 0$
- Only do it if this doesn't introduce bias into your sample
 - In general, only do it if $x_i > 0$ and $y_i > 0$ for almost all i
 - Adding 1 or 0.1, or 100 is not a valid fix**

22

An example

- A regression of ln(wages) on: ln(Age), race (black=1) and ln(IQ) (IQ is the percentile)
- For an increase in one percent in age, we expect wages to change by $\hat{\beta}_{age}$ percent
- On average, we expect wages to be higher/lower for blacks by $\hat{\beta}_{black}$ percent than for non-blacks
Black
- For an increase in one percent in the IQ percentile (that is, a percent change in percentage points), we expect wages to change by $\hat{\beta}_{IQ}$ percent

23

Simulations!

```
loglog=lm(log(wage) ~ log(IQ_Percentile) + log(age) + black, data = wage2)
summary(loglog)
stargazer(loglog, title="Log-Level", align=TRUE,
  type="latex", omit.table.layout="la",
  out="Lectures/tables/loglog.tex",
  covariate.labels=c("ln(IQ (percentile))", "ln(Age)", "Black(=1)"),
  digits=2, digits.extra=1, no.space=T, colnames=F,
  dep.var.caption="", dep.var.labels="ln(Wage)",
  column.sep.width="0pt", header=F,
  omit.stat=c("adj.rsq", "rsq", "F", "ser"))
```

24

```
Call:
lm(formula = log(wage) ~ log(IQ_Percentile) + log(age) + black,
    data = wage2)

Residuals:
    Min       1Q   Median       3Q      Max
-1.98259 -0.25865  0.01121  0.28098  1.30397

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  4.3983406  0.5531443   7.923 6.56e-15 ***
log(IQ_Percentile) -0.0009437  0.0013559  -0.018  0.985
log(age)      0.6929047  0.1412305   4.906 1.10e-06 ***
black        -0.2850449  0.0396476  -7.189 1.33e-12 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.4051 on 931 degrees of freedom
Multiple R-squared:  0.07774    Adjusted R-squared:  0.07477
F-statistic: 26.16 on 3 and 931 DF,  p-value: 2.994e-16
```

25

Log-Level

	ln(Wage)
ln(IQ (percentile))	-0.001 (0.05)
ln(Age)	0.69*** <i>0.69%</i> (0.14)
Black(=1)	-0.29*** <i>-29%</i> (0.04)
Constant	4.40*** (0.56)
Observations	935
Note:	*p<0.1; **p<0.05; ***p<0.01

26

Beyond the basic of OLS

A few things that don't get enough attention

How to interpret coefficients/regression table

Leverage

The perils of p-hacking

What if your outcome is a dummy?

Ordinal/Categorical data

Error structure

Heteroskedasticity

Cluster standard errors

Statistical power

Randomizing at the Unit of Analysis

Cluster Randomized Experiments

27

Leverage

- Remember that

$$\hat{\beta} = \frac{\text{cov}(x, y)}{v(x)} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

- We can rewrite as:

$$\hat{\beta} = \frac{(x_1 - \bar{x})(y_1 - \bar{y})}{(x_1 - \bar{x})^2} + \frac{\sum_{i=2}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=2}^n (x_i - \bar{x})^2}$$

- If $x_1 = \bar{x}$, then

$$\hat{\beta} = \frac{\sum_{i=2}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=2}^n (x_i - \bar{x})^2}$$

- The first observation doesn't affect the outcome

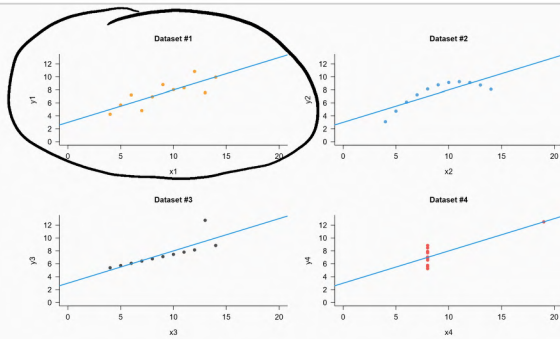
28

Leverage: Big Picture

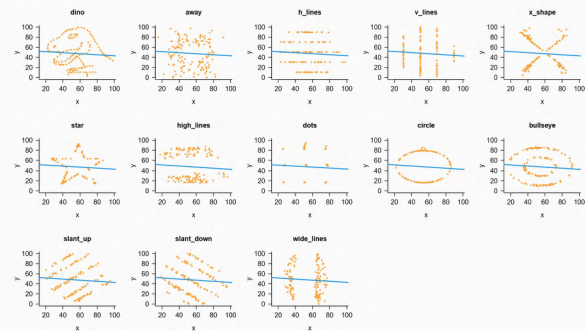
- That was an extreme case ($x_i = \bar{x}$) but generally speaking:
- The farther an observation is from \bar{x} , the more it affects the OLS estimator
- This is called "leverage"

- See a recent discussion on Twitter of economist arguing about this <https://twitter.com/arindube/status/1279919438419165184?s=20>

29



30



31

Beyond the basic of OLS

A few things that don't get enough attention

How to interpret coefficients/regression table

Leverage

The perils of p-hacking

What if your outcome is a dummy?

Ordinal/Categorical data

Error structure

Heteroskedasticity

Cluster standard errors

Statistical power

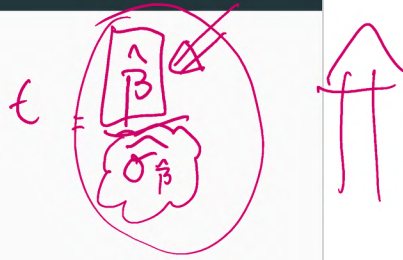
Randomizing at the Unit of Analysis

Cluster Randomized Experiments

32

The perils of p-hacking

<https://xkcd.com/882/>



33

Beyond the basic of OLS

A few things that don't get enough attention

How to interpret coefficients/regression table

Leverage

The perils of p-hacking

What if your outcome is a dummy?

Ordinal/Categorical data

Error structure

Heteroskedasticity

Cluster standard errors

Statistical power

Randomizing at the Unit of Analysis

Cluster Randomized Experiments

34

What if your outcome is a dummy?

- All we have talked about still holds
- Logit/Probit have very strong assumptions (the shape of the error term)
- Regression is more robust in general

35

Beyond the basic of OLS

A few things that don't get enough attention

Error structure

Statistical power

40

Beyond the basic of OLS

A few things that don't get enough attention

Error structure

Statistical power

41

Variance of OLS estimators

The correct variance estimation procedure is given by the structure of the data

- It is very unlikely that all observations in a dataset are unrelated, but drawn from identical distributions (**homoskedasticity**) i.i.d. $\epsilon \text{ es } \epsilon.i.d.$
- For instance, the variance of income is often greater in families belonging to top deciles than among poorer families (**heteroskedasticity**)
- Some phenomena do not affect observations individually, but they do affect groups of observations uniformly within each group (**clustered data**)

Beyond the basic of OLS

A few things that don't get enough attention

How to interpret coefficients/regression table

Leverage

The perils of p-hacking

What if your outcome is a dummy?

Ordinal/Categorical data

Error structure

Heteroskedasticity

Cluster standard errors

Statistical power

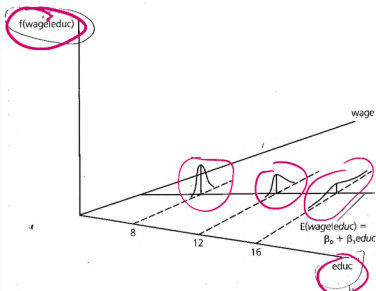
Randomizing at the Unit of Analysis

Cluster Randomized Experiments

43

OLS inference is generally faulty in the presence of heteroskedasticity

Figure 2.9
Var (wage|educ) increasing with educ.



Heteroskedasticity

- Assume

$$\text{Var}(u_i|x_i) = \sigma_i^2$$

$\text{VAR}(u_i|x_i) = \sigma^2$
HOMOSKEDASTICIDAD

- Fortunately, OLS is still useful ($\hat{\beta}$ still consistent/unbiased)
- Note that errors are still independent from each other
- The variance of our estimator, $\hat{\beta}_1$ equals:

$$\text{Var}(\hat{\beta}_1) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = (X'X)^{-1} X' V(u_i|x_i) X (X'X)^{-1}$$

- When $\sigma_i^2 = \sigma^2$ for all i , this formula reduces to the usual form,
 $\frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = \sigma^2 (X'X)^{-1}$

Robust standard errors

- A valid estimator of $\text{Var}(\hat{\beta}_1)$ for heteroskedasticity of any form (including homoskedasticity) is

$$\text{Var}(\hat{\beta}_1) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \hat{u}_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = (X'X)^{-1} X' \left(\sum_{i=1}^n x_i x_i' \hat{u}_i^2 \right) X (X'X)^{-1}$$

which is easily computed from the data after the OLS regression

- As a rule, you should always use "robust standard errors"

Simulations!

```
library(sandwich)
alpha=0 #intercept
Reps=1000 #how many simulations?
Nobs=100 #number of obs
SequenceBetas=seq(0.1,0.1) #lets do different betas
FractionSignificant=NULL #fraction significant 5% level
FractionSignificant_robust=NULL #fraction significant 5% level when using robust
betaVector=NULL #mean estimator
betaVector_robust=NULL #mean estimator robust
```

Simulations!

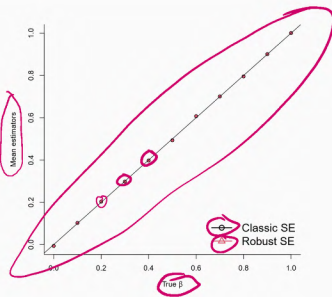
```

for(beta in SequenceBetas){
  #save the outcomes from the simulations
  beta_estimate=rep(NA,Reps)
  beta_pvalue=rep(NA,Reps)
  beta_estimate_robust=rep(NA,Reps)
  beta_pvalue_robust=rep(NA,Reps)
  #generate some x data
  xas_matrix(runif(Nobs,-5.5))
  for(r in 1:Reps){
    #use the DGP to generate outcome data with heteroskedasticity
    Y=alpha+beta*x+sigma*norm(Nobs, sd=1)*x
    OLS=lm(Y~X) #estimate OLS
    ResultsOLS=summary(OLS)$coef #save results from OLS table
    beta_estimate[r]=ResultsOLS[2,1]
    beta_pvalue[r]=ResultsOLS[2,4]
    #Results from robust OLS: HCl yields same results as stata
    ResultsRobust=coefest(OLS, vcov = vcovHC(OLS, type = "HCl"))
    beta_estimate_robust[r]=ResultsRobust[2,1]
    beta_pvalue_robust[r]=ResultsRobust[2,4]
  }
  #Save the results for the given value of beta
  FractionSignificant=(FractionSignificant, mean(beta_pvalue<0.05))
  FractionSignificant_robust=(FractionSignificant_robust, mean(beta_pvalue_robust<0.05))
  betaVector=(betaVector, mean(beta_estimate))
  betaVector_robust=(betaVector_robust, mean(beta_estimate_robust))
}

```

48

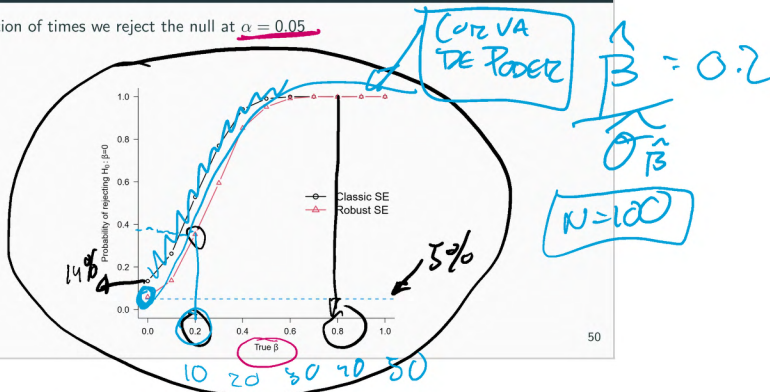
No bias



49

Power Curve – Incorrect type-I error from classic OLS, correct from robust SE

Proportion of times we reject the null at $\alpha = 0.05$



50

Beyond the basic of OLS

A few things that don't get enough attention

How to interpret coefficients/regression table

Leverage

The perils of p-hacking

What if your outcome is a dummy?

Ordinal/Categorical data

Error structure

Heteroskedasticity

Cluster standard errors

Statistical power

Randomizing at the Unit of Analysis

Cluster Randomized Experiments

51

Clustered data

- But what if errors are not independent?
- Maybe observations between units in a group are related to each other
 - Imagine you randomly assign a treatment at the school level (e.g., extra resources)
 - The **unobservables** of kids belonging to the same school are correlated (e.g., teacher quality, recess routines)
 - The **unobservables** of kids in different school are unlikely to be correlated
- Then independence of errors across observations is violated
- But maybe independence holds across schools, just not within schools

ESTADOS
- COLOMBIA

Simulations!

```
Classes=50 #number of classes or schools
StudentsPerClass=10 #number of obs per schools
Reps=1000 #repetitions
SequenceBetas=seq(0,1,0.1) #try different betas (treatment effects)
alpha=0 #intercept
FractionSignificant=NULL #fraction significant 5% level
FractionSignificant_robust=NULL #fraction significant 5% level when using robust
betaVector=NULL #mean estimator
betaVector_robust=NULL #mean estimator robust
```

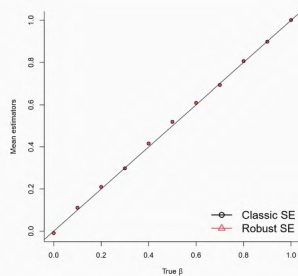
53

Simulations!

```
for(beta in SequenceBetas){
  #save the outcomes from the simulations
  beta_estimate=rep(NA,Reps)
  beta_pvalue=rep(NA,Reps)
  beta_estimate_robust=rep(NA,Reps)
  beta_pvalue_robust=rep(NA,Reps)
  X=as.matrix(runif(StudentsPerClass*Classes,-5,5)) #generate some x data
  for(r in 1:Reps){
    Schoks_Cluster=rep(1:norm(Classes),each=StudentsPerClass)
    Schoks_Individual=norm(StudentsPerClass*Classes,sd=1)
    Y=alpha+beta*X+Schoks_Cluster+Schoks_Individual
    OLS=lm(Y~X) #estimate OLS
    ResultsOLS=summary(OLS)$coef
    beta_estimate[r]=ResultsOLS[2,1]
    beta_pvalue[r]=ResultsOLS[2,4]
    #Results from robust OLS: HCL yields same results as stata
    ResultsRobust=covfct(OLS,vcov=vcovHC(OLS, type="HCL"))
    beta_estimate_robust[r]=ResultsRobust[2,1]
    beta_pvalue_robust[r]=ResultsRobust[2,4]
  }
  #Save the results for the given value of beta
  FractionSignificant=c(FractionSignificant,mean(beta_pvalue<0.05))
  FractionSignificant_robust=c(FractionSignificant_robust,mean(beta_pvalue_robust<0.05))
  betaVector=c(betaVector,mean(beta_estimate))
  betaVector_robust=c(betaVector_robust,mean(beta_estimate_robust))
}
```

54

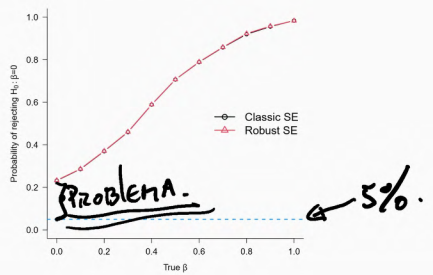
No bias



55

Power Curve – Incorrect type-I error from classic OLS and from robust SE

Proportion of times we reject the null at $\alpha = 0.05$



Cluster robust standard errors

- Both classic OLS and robust SE overreject (i.e., they reject the null when its true more times than we thought at a given level)
- We need to allow for arbitrary correlation within group
- Instead of summing over each individual, we first sum over groups
- I'll use matrix notation as it's easier for me to explain by stacking the data

$$\sigma^2 = \frac{\sum_{i=1}^n u_i^2}{n-1}$$

Clustered data

- Let's stack the observations by cluster

$$y_g = x_g \beta + u_g \quad \text{gy escuelas.}$$

- The OLS estimator of β is:

$$\hat{\beta} = [X'X]^{-1} X'y$$

- The variance is given by:

$$\text{Var}(\hat{\beta}) = E[[X'X]^{-1} X' \Omega X [X'X]^{-1}]$$

$\Omega = V(u)$

Clustered data

With this in mind, we can now write the variance-covariance matrix for clustered data

$$\text{Var}(\hat{\beta}) = [X'X]^{-1} \left[\sum_{g=1}^G X'_g \hat{u}_g \hat{u}'_g X_g \right] [X'X]^{-1}$$

where \hat{u}_g are residuals from the stacked regression

Formula Heteroskedasticidad para el grupo g

- In STATA: vce(cluster clustervar)
- In R use lfe package

$$X'_g u_g \hat{u}'_g X_g \Rightarrow \text{Implicitamente una sumatoria}$$

Simulations!

```
library(lfe)
Classes=50 #number of classes or schools
StudentsPerClass=5 #number of obs per schools
Reps=1000 #repetitions
SequenceBetas=seq(0.1,0.1) #try different betas (treatment effects)
alpha=0 #intercept
FractionSignificant=NULL #fraction significant 5% level
FractionSignificant_robust=NULL #fraction significant 5% level when using robust
FractionSignificant_cluster=NULL #fraction significant 5% level when using cluster
```

60

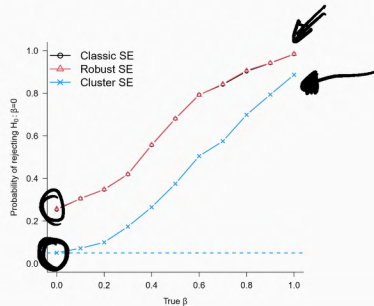
Simulations!

```
for(beta in SequenceBetas){
  #save the outcomes from the simulations
  beta_pvalue=rep(NA,Reps)
  beta_pvalue_robust=rep(NA,Reps)
  beta_pvalue_cluster=rep(NA,Reps)
  ClusterIndicator=rep(1,Classes,each=StudentsPerClass)
  TreatmentClassLevel=sample(0,1,Classes,replace=T)
  TreatmentIndividual=rep(TreatmentClassLevel,each=StudentsPerClass)
  for(r in 1:Reps){
    Schoks_Cluster=rep(rnorm(Classes),each=StudentsPerClass)
    Schoks_Individual=rnorm(StudentsPerClass*Classes,sd=1)
    #alpha+beta*TreatmentIndividual+Schoks_Cluster+Schoks_Individual
    OLS=lm(Y~TreatmentIndividual | 0 | 0 | ClusterIndicator) #estimate OLS
    beta_pvalue[r]=OLS$pval[2]
    #Results from robust SE
    beta_pvalue_robust[r]=OLS$r.pval[2]
    #Results from cluster SE
    beta_pvalue_cluster[r]=OLS$c.pval[2]
  }
  #Save the results for the given value of beta
  FractionSignificant=(FractionSignificant,mean(beta_pvalue<0.05))
  FractionSignificant_robust=(FractionSignificant_robust,mean(beta_pvalue_robust<0.05))
  FractionSignificant_cluster=(FractionSignificant_cluster,mean(beta_pvalue_cluster<0.05))
}
```

61

Power Curve)

Proportion of times we reject the null at $\alpha = 0.05$



62

The importance of knowing your data

- In real world you should never go with the “independent and identically distributed” (i.e., homoskedasticity) case. Life is not that simple.
- You need to know your data in order to choose the correct error structure and then infer the required SE calculation
- At a minimum, use robust standard errors
- If you have aggregate variables, like class size, you need to consider clustering at that level

When to cluster?

- Case 1: If sampling follows a two stage process where in the first stage, a subset of clusters were sampled randomly from a population of clusters, and in the second stage, units were sampled randomly from the sampled clusters
- Case 2: When clusters of units, rather than units, are assigned to a treatment

When to cluster?

- The results on cluster SE

$$\text{Var}(\hat{\beta}) = [X'X]^{-1} \left[\sum_{g=1}^G x_g' \hat{u}_g \hat{u}_g' x_g \right] [X'X]^{-1}$$

relies on "asymptotic results" based on the number of clusters (G) — not on the total sample size N

- Can only use cluster SE if number of clusters is "large" (usually over $\sim 40 - 50$)
- If number of clusters is small consider:
 - Collapsing the data at the "cluster" level
 - Wild bootstrap
 - Randomization inference (if you have an experiment)

When to cluster?

- Two good reads on clustering:

⇒ • Cameron, A.C. and Miller, D.L., 2015. A practitioner's guide to cluster-robust inference. *Journal of human resources*.
<http://jhr.oup.com/content/50/2/317.refs>

⇒ • Abadie, A., Athey, S., Imbens, G.W. and Wooldridge, J., 2017. When should you adjust standard errors for clustering? (No. w24003). National Bureau of Economic Research. <https://www.nber.org/papers/w24003>

Beyond the basic of OLS

A few things that don't get enough attention

Error structure

Statistical power

Beyond the basic of OLS

A few things that don't get enough attention

Error structure

Statistical power

68

Introduction

- In a simple experiment the average treatment effect is the difference in sample means between the treatment and the control group
- This is the OLS coefficient of β in the regression

$$Y_i = \alpha + \beta T_i + \varepsilon_i$$

69

Regression analysis of OLS

And

$$X'X = \sigma^2 \begin{pmatrix} \frac{1}{p} & 1 \\ 1 & 1 \end{pmatrix}$$
$$N_T = pN$$
$$N_C = (1-p)N$$
$$(X'X)^{-1} = \frac{1}{N(1-p)} \begin{pmatrix} 1 & -1 \\ -1 & \frac{1}{p} \end{pmatrix}$$
$$V \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = \sigma^2 (X'X)^{-1} = \sigma^2 \frac{1}{N(1-p)} \begin{pmatrix} 1 & -1 \\ -1 & \frac{1}{p} \end{pmatrix}$$
$$V(\hat{\beta}) = \sigma^2 \frac{1}{Np(1-p)}$$

Statistical power

How many observations are enough?

71

Statistical power

How many observations are enough?

Definition

The **power of the design** is the probability that, for a given effect size and a given statistical significance level, we will be able to reject the hypothesis of zero effect

$$P(\text{Reject } H_0 \mid \beta = \beta_0)$$

71

Statistical power

- Is the unit of treatment the same as the unit of analysis? Or, is the treatment to be administered to a 'cluster' of units?

72

Statistical power

- Is the unit of treatment the same as the unit of analysis? Or, is the treatment to be administered to a 'cluster' of units?
- Examples of individual randomizations:
 - Individuals who are given mobile phones to induce them to use an m-banking platform
 - Farmers individually provided with improved agricultural inputs
 - Students admitted to an elite school by a lottery process

72

Beyond the basic of OLS

A few things that don't get enough attention

How to interpret coefficients/regression table

Leverage

The perils of p-hacking

What if your outcome is a dummy?

Ordinal/Categorical data

Error structure

Heteroskedasticity

Cluster standard errors

Statistical power

Randomizing at the Unit of Analysis

Cluster Randomized Experiments

73

Randomizing at the Unit of Analysis

- The estimate of treatment effect is $\hat{\beta}$ in the regression

$$Y_i = \alpha + \beta T_i + \varepsilon_i$$

- The mean of $\hat{\beta}$ is β (the true effect)
- The variance of $\hat{\beta}$ is $V(\hat{\beta}) = \frac{\sigma^2}{p(1-p)N}$
- σ^2 is the variance of the outcome (Y_i)
- p is the proportion of treated units
- N is the number of observations

74

Randomizing at the Unit of Analysis

- We are generally interested in testing the null hypothesis (H_0) that the effect of the program is equal to zero against the alternative that it is not
- The **significance level**, or size, of a test represents the probability of a type I error, i.e., the probability we reject the hypothesis when it is in fact true
- The **power of the test** the probability that we reject H_0 when it is in fact false

75

Randomizing at the Unit of Analysis

- We are generally interested in testing the null hypothesis (H_0) that the effect of the program is equal to zero against the alternative that it is not
- The **significance level**, or size, of a test represents the probability of a type I error, i.e., the probability we reject the hypothesis when it is in fact true
- The **power of the test** the probability that we reject H_0 when it is in fact false

We will constantly use the fact that:

$$\hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{p(1-p)N}\right)$$

75

Randomizing at the Unit of Analysis

- We are generally interested in testing the null hypothesis (H_0) that the effect of the program is equal to zero against the alternative that it is not
- The **significance level**, or size, of a test represents the probability of a type I error, i.e., the probability we reject the hypothesis when it is in fact true
- The **power of the test** the probability that we reject H_0 when it is in fact false

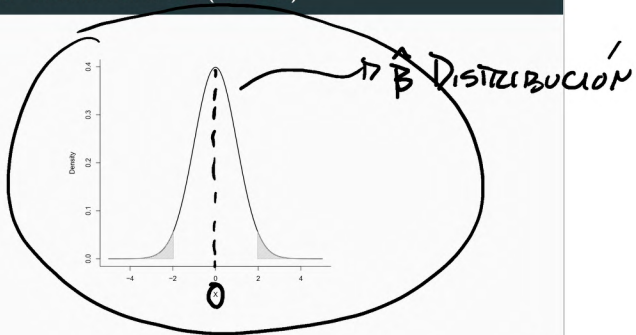
We will constantly use the fact that:

$$\hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{p(1-p)N}\right)$$

We often normalize the outcome and present results in terms of SD (so $\sigma^2 = 1$).

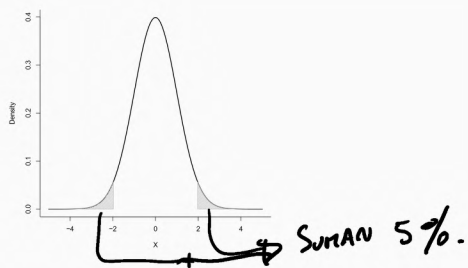
75

Significance level - Assume null is true (no effect)



76

Significance level - Assume null is true (no effect)



Gray area is the probability we reject the null when it is true

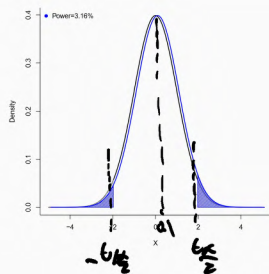
76

Power when the effect is β_1

For a true effect size β this is the fraction of the area under this curve that falls to the right of the critical value $t_{\alpha/2}$

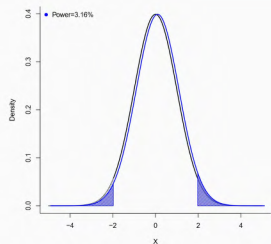
77

Power when the effect is $\beta = 0.1$



78

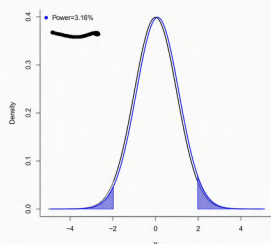
Power when the effect is $\beta = 0.1$



Blue area is the probability we reject the null when β is 0.1

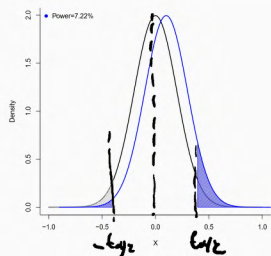
78

Power when $\beta = 0.1$, $N = 4$, and $p = 0.5$



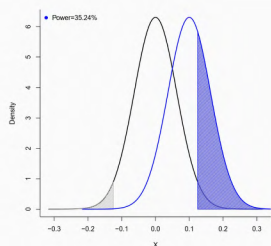
79

Power when $\beta = 0.1$, $N = 100$, and $p = 0.5$



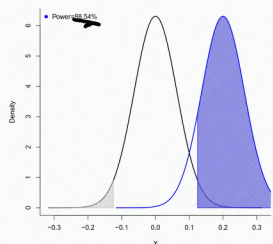
80

Power when $\beta = 0.1$, $N = 1,000$, and $p = 0.5$



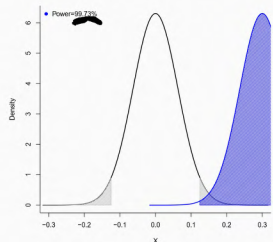
81

Power when $\beta = 0.2$, $N = 1,000$, and $\rho = 0.5$



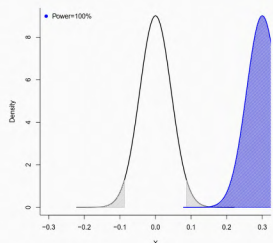
Blue area is the probability we reject the null when β is 0.2

Power when the effect is $\beta = 0.3$, $N = 1,000$, and $\rho = 0.5$



Blue area is the probability we reject the null when β is 0.3

Power when the effect is $\beta = 0.3$, $N = 1,000$, $\rho = 0.5$, and $\sigma = 0.7$



Blue area is the probability we reject the null when β is 0.3

Statistical power and clusters

- All these quantities we just looked at are related

- To achieve a power $(1-\alpha)$ it must therefore be that

$$\beta > (t_{\frac{\alpha}{2}} + t_{1-\alpha}) \sigma_{\hat{\beta}}$$

\downarrow SE($\hat{\beta}$)
 \downarrow
 $\hookrightarrow \alpha = 5\%$
1.96

Minimum detectable effect

- The minimum detectable effect size for a given power (κ), significance level (α), sample size (N), and portion of subjects allocated to treatment group (p) is given by

$$MDE = (t_{\frac{\alpha}{2}} + t_{1-\kappa}) \sqrt{\frac{\sigma^2}{p(1-p)N}}$$

$SE(\hat{\beta})$

$\alpha = 0.1\%$ → $80\% = \kappa$

$$SE(\hat{\beta}) = \frac{\sigma^2(N)}{N_T N_C}$$

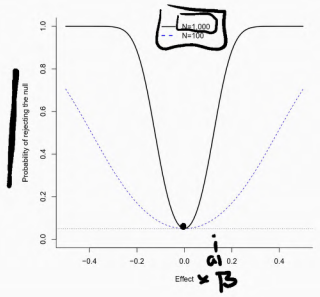
$$\frac{\partial SE(\hat{\beta})}{\partial N_T} = 0 \rightarrow \frac{-(1)N_C + (-1)N_T}{(N_T N_C)^2}$$

$N = N_T + N_C$
 $\partial N_T / \partial N_T$

Randomizing at the Unit of Analysis

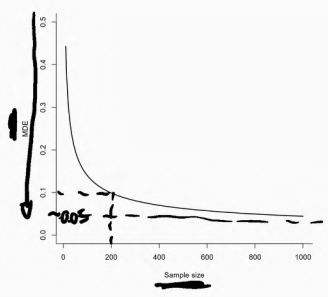
- The standard is to set $\kappa = 0.8$ or $\kappa = 0.9$
- The standard is to set $\alpha = 0.05$ or $\alpha = 0.1$
- The variance of outcomes σ^2 is typically the raw variance of the dependent variable you intend to use
- The sample size N is the number of observations in the study (you can change this)
- The fraction of the sample treated is p (you can change this)

Effect vs Power



$\sigma = 1$
 $p = 0.5$

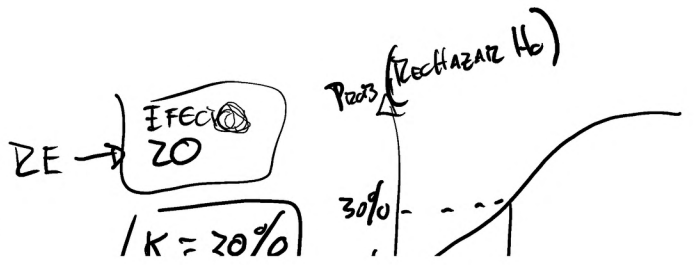
Sample size vs MDE



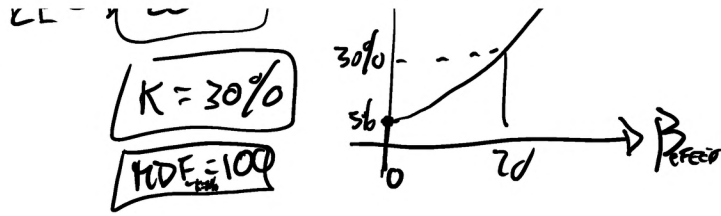
$p = 0.5$
 $\sigma = 1$
 $\alpha = 5\%$
 $\kappa = 90\%$

How should you think about the MDE?

- What is the treatment effect below which it is pointless to implement the program



- What is the treatment effect below which it is pointless to implement the program and/or study its effect?
- If sample size is too small, you're likely to end up with an insignificant result for something that actually matters



$HDE = 0.0001$
 $\hat{\beta} = 0.0001 \text{ (P-val } < 0.01)$

90

Beyond the basic of OLS

- A few things that don't get enough attention
 - How to interpret coefficients/regression table
 - Leverage
 - The perils of p-hacking
 - What if your outcome is a dummy?
 - Ordinal/Categorical data
- Error structure
 - Heteroskedasticity
 - Cluster standard errors
- Statistical power
 - Randomizing at the Unit of Analysis
 - Cluster Randomized Experiments

91

Cluster Randomized Experiments

- Is the unit of treatment the same as the unit of analysis? Or, **is the treatment to be administered to a 'cluster' of units?**

92

Cluster Randomized Experiments

- Is the unit of treatment the same as the unit of analysis? Or, **is the treatment to be administered to a 'cluster' of units?**
- Examples of clustered randomizations:
 - Changing the business practices at a firm level and studying the impact on individual employees
 - Providing schools with new textbooks and studying the effect on individual student performance
 - Offering a new financial service to all residents in a village and studying the impact on micro enterprise outcomes
- In a clustered randomization the power of the study is coming partly from the number of individuals in the study, and partly from the number of clusters in the study

92

Cluster Randomized Experiments

- The estimate of treatment effect is $\hat{\beta}$ in the regression

$$Y_{ij} = \alpha + \beta T_j + \omega_j + \varepsilon_{ij}$$

- σ^2 is the variance of the outcome (ε_{ij})
- τ^2 is the variance of the outcome (ω_j)
- p is the proportion of treated units
- n is the number of observations in each cluster
- J is the number of clusters

The variance of $\hat{\beta}$ is $\sigma_{\hat{\beta}}^2 = \frac{n\tau^2 + \sigma^2}{p(1-p)J}$

SE Cluster? $SE(\hat{\beta}) = \frac{\sigma^2}{p(1-p)N}$

93

Cluster Randomized Experiments

- Often, expressed using the intra-cluster correlation (ICC) $\equiv \frac{\tau^2}{\tau^2 + \sigma^2} = \rho$
- The variance of $\hat{\beta}$ is $V(\hat{\beta}) = \sigma^2 \left(\frac{1-p}{p(1-p)J} \right)$ (comes from the cluster SE formula we saw)

- The ICC can be obtained using *loneway* in stata

94

RHO

SI $n \rightarrow \infty \Rightarrow V(\hat{\beta}) = \frac{p \sigma^2}{p(1-p)J}$

SI $J \rightarrow \infty \Rightarrow V(\hat{\beta}) = 0$ [sin impossible]

Minimum detectable effect

- The **minimum detectable effect** is given by

$$MDE = (t_{\alpha/2} + t_{1-\beta}) \sigma \sqrt{\frac{\rho + \frac{(1-\rho)}{n}}{p(1-p)J}}$$

SE($\hat{\beta}$)

$n=10, s=200 \rightarrow MDE$

$n=5, s=400 \rightarrow MDE$

95

Power Calculations Rules of Thumb

- For an individual-level experiment, 200-300 observations will typically be sufficient to detect a reasonable effect size
- For a clustered experiment, a low ICC (0.1) would need 50-100 clusters and > 5 observations per cluster to detect a moderate effect. As the ICC gets larger, the number of **clusters** has to go up
- For **very** complicated research designs, you can always use simulations to get the power of the design

96