

Political Science 209 - Fall 2018

Measurement

Florian Hollenbach

25th September 2018

- A sample is a small share of the population in that we are interested in

Survey Sampling

- A sample is a small share of the population in that we are interested in
- How do we draw samples in such a way that polls accurately reflect what is going to happen?
- How to construct samples that will represent the population?

- Example: We want to know the voting intentions of Texans (or Americans)
- We can hardly ask all eligible voters about their intention

Survey Sampling

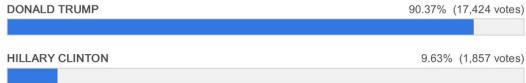
- Example: We want to know the voting intentions of Texans (or Americans)
- We can hardly ask all eligible voters about their intention
- We take a *sample*

Survey Sampling

- The size of the sample is less important than its composition



DRUDGE POLL WHO WON THE FIRST PRESIDENTIAL DEBATE?



Total Votes: 19,281

Suka Kongsi 607

Tweet

- Mail questionnaire to 10 million people
- Addresses came from phone books and club memberships
- Problems?

Literary Digest Sample

- Mail questionnaire to 10 million people
- Addresses came from phone books and club memberships
- Problems?
- Biased *sample*

Quota Sampling

- Sample certain groups until quota is filled
- Does not mean unobservables are representative

Simple Random Sampling

- Think of all voters sitting in a box, survey firm randomly draws voters
- Random draws without replacement give us an unbiased estimate of the population
- Everybody has the same chance of being in the sample

Simple Random Sampling

- Pre-determined number of units are randomly selected from population
- Sample will be representative of population on observed and unobserved characteristics

Simple Random Sampling

- Not every single sample will be exactly representative
- If we were to take a lot of random samples (say 1000 samples of 1000 respondents), on average the samples would be representative

Simple Random Sampling

- Each single sample can be off and different
- Polls are associated with uncertainty

Ted Cruz leads Beto O'Rourke 54 to 45, new poll says

The new Quinnipiac University poll surveyed likely voters instead of registered voters like it did in past iterations.

BY **PATRICK SVITEK** SEPT. 18, 2018 11 AM

Simple Random Sampling

- Each single sample can be off and different
- Polls are associated with uncertainty

Ted Cruz leads Beto O'Rourke 54 to 45, new poll says

The new Quinnipiac University poll surveyed likely voters instead of registered voters like it did in past iterations.

BY [PATRICK SVITEK](#) SEPT. 18, 2018 11 AM

Beto O'Rourke leads Ted Cruz by 2 among likely voters in U.S. Senate race, new poll finds

O'Rourke has been closing the gap over the last several months, but this is the first poll that puts him ahead of Cruz.




BY [KATHRYN LUNDSTROM](#) SEPT. 19, 2018 8 AM

Random Sampling is hard

- How to create sampling frame?
- Random digit dialing? Walking to random houses?
- Multi-stage cluster sampling

Non-response bias

- Unit non-response bias:

 **Nate Cohn** 
@Nate_Cohn [Follow](#) 

We've called thousands of cell phones in VA07 today and we don't have a single new 18-29 year old respondent. Meanwhile, we've got a disproportionate 8 in Colo 6 already, all supporting the Dem

| Age | | DEM. | REP. | UND. |
|--------------|-----------------------|------|------|------|
| 18 to 29 | n = 8 / 14% of voters | 84% | — | 16% |
| 30 to 44 | 13 / 26% | 74% | 22% | 4% |
| 45 to 64 | 16 / 36% | 24% | 64% | 12% |
| 65 and older | 11 / 24% | 36% | 64% | — |

8:13 PM - 12 Sep 2018

Non-response bias

- Item non-response bias: *What was the last crime you committed?*
- Sensitive questions: non-response, social desirability bias
Turnout, racial prejudice, corruption

Why could this be a problem in the Afghanistan example?



Examples of Problems of Running Surveys in Afghanistan

- unit non-response bias:

Examples of Problems of Running Surveys in Afghanistan

- unit non-response bias:

some citizens might not want to answer the door

Examples of Problems of Running Surveys in Afghanistan

- unit non-response bias:

some citizens might not want to answer the door

- item non-response bias:

Examples of Problems of Running Surveys in Afghanistan

- unit non-response bias:

some citizens might not want to answer the door

- item non-response bias:

Taliban supporters may be less likely to answer questions about Taliban

Examples of Problems of Running Surveys in Afghanistan

- unit non-response bias:

some citizens might not want to answer the door

- item non-response bias:

Taliban supporters may be less likely to answer questions about Taliban

- social desirability bias:

Examples of Problems of Running Surveys in Afghanistan

- unit non-response bias:

some citizens might not want to answer the door

- item non-response bias:

Taliban supporters may be less likely to answer questions about Taliban

- social desirability bias:

Taliban supporters may not want to admit to supporting Taliban

List Experiments

- list of groups respondent might support
- Asked to name number of groups they support
- Treated subject with controversial group, control group without controversial group

Strategies to Ask Sensitive Questions

List Experiments - Control

I'm going to read you a list with the names of different groups and individuals on it. After I read the entire list, I'd like you to tell me how many of these groups and individuals you broadly support, meaning that you generally agree with the goals and policies of the group or individual. Please don't tell me which ones you generally agree with; only tell me how many groups or individuals you broadly support.

Groups: Karzai Government; National Solidarity Program; Local Farmers

Strategies to Ask Sensitive Questions

List Experiments - Treated

I'm going to read you a list with the names of different groups and individuals on it. After I read the entire list, I'd like you to tell me how many of these groups and individuals you broadly support, meaning that you generally agree with the goals and policies of the group or individual. Please don't tell me which ones you generally agree with; only tell me how many groups or individuals you broadly support.

Groups: Karzai Government; National Solidarity Program; Local Farmers;
ISAF (Taliban)

List Experiments

- Average difference between Treated and Control group is the estimated percentage of people who support controversial group

Summarizing Bivariate Relationships

- Bivariate relationships are associations between **two** variables
- Example: treatment of Spanish confederates (X or T) and exclusionary attitudes (Y)

Simple Summaries of Bivariate Relationships

If X (independent variable) is categorical:

- Comparison of means
- boxplots

Simple Summaries of Bivariate Relationships

If both X (independent variable) and Y (dependent variable) are continuous:

Simple Summaries of Bivariate Relationships

If both X (independent variable) and Y (dependent variable) are continuous:

- Scatterplots

Simple Summaries of Bivariate Relationships

If both X (independent variable) and Y (dependent variable) are continuous:

- Scatterplots
- Correlation

Simple Summaries of Bivariate Relationships: Scatterplot

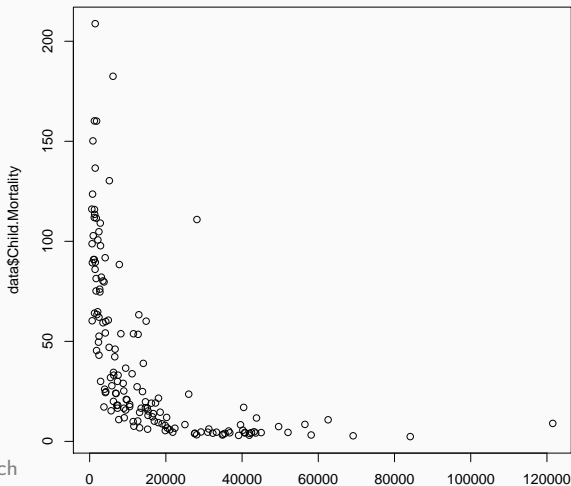
- Direct graphical comparison of two variables
- Use `plot(y,x)` in *R*

Scatterplot

```
data <- read.csv("bivariate_data.csv")  
data <- subset(data, year == 2010)  
plot(data$GDP,data$Child.Mortality)
```

Scatterplot

That looks weird, no? What do we do with skewed variables?



What do we do with skewed variables?

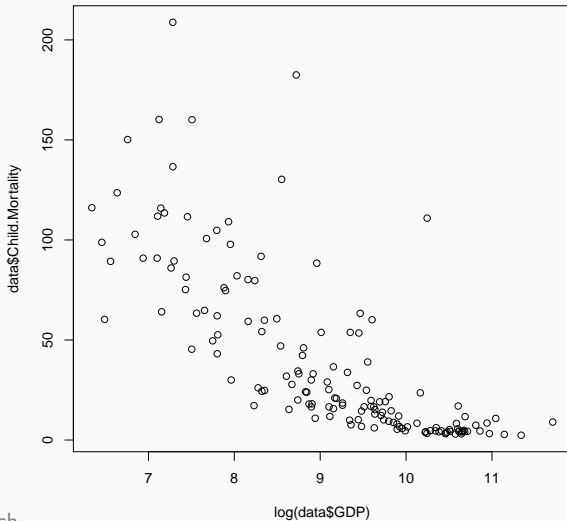
What do we do with skewed variables?

- When variable have a small number of observations with extremely large or small positive values, we often take the natural log
- The natural logarithm is the logarithm with base e , which is a mathematical constant approximately equal to 2.7182 (inverse of e^y)

Scatterplot

```
plot(log(data$GDP),data$Child.Mortality)
```

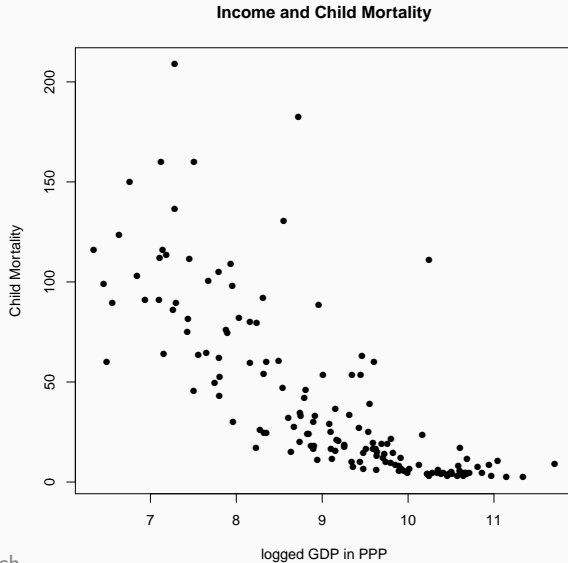
Scatterplot



Scatterplot

```
pdf("~/Documents/GitHub/Polisci209_2018/slides/week5/scatter.pdf")
plot(log(data$GDP), data$Child.Mortality, pch = 16, col = "black",
      xlab = "logged GDP in PPP", ylab = "Child Mortality", main = "Income and Child Mortality")
dev.off()
```

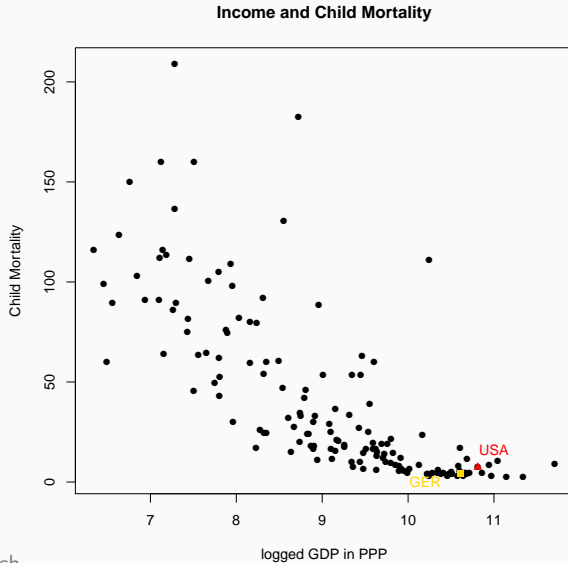
Scatterplot



Scatterplot – more fun

```
## add special points for USA and Germany
pdf("~/Documents/GitHub/Polisci209_2018/slides/week5/scatter_points.pdf")
plot(log(data$GDP), data$Child.Mortality, pch = 16, col = "black",
     xlab = "logged GDP in PPP", ylab = "Child Mortality", main = "Income and Child Mortality")
points(log(data$GDP[data$Country.code == "USA"]), data$Child.Mortality[data$Country.code == "USA"], pch = 17, col = "red") ##USA
text(11, 16, "USA", col = "red")
points(log(data$GDP[data$Country.code == "DEU"]), data$Child.Mortality[data$Country.code == "DEU"], pch = 15, col = "gold") ##USA
text(10.2, 0, "GER", col = "gold")
dev.off()
```

Scatterplot



Bivariate Relationships in Numbers

- How can we quantify the relationship between two continuous variables?

Bivariate Relationships in Numbers

- How can we quantify the relationship between two continuous variables?
- Correlation: the most used measure of bivariate relationships
- Correlation measures how two variables move together relative to their respective means.

Calculating correlations – Step 1: Standardizing a variable

- By standardizing we bring all variables on the same scale
- The resulting mean will be zero, the standard deviation will be one
- We standardize by subtracting a variables mean and dividing by the standard deviation

Calculating correlations – Step 1: Standardizing a variable

- Standardized variables are also called z-scores:

$$Z_i = \frac{x_i - \bar{x}(\text{mean of } x)}{sd_x(\text{standard deviation of } x)}$$

- The z-score is independent of the scale of the variable or shifts in the variable

Calculating correlations – Step 1: Standardizing a variable

- Standardized variables are also called z-scores:

$$Z_i = \frac{x_i - \bar{x}(\text{mean of } x)}{sd_x(\text{standard deviation of } x)}$$

- The z-score is independent of the scale of the variable or shifts in the variable
- This means GDP and (GDP*100 + 10000) will have the exact same z-scores

Calculating correlations – Step 2: Calculating the correlation

$$\text{Correlation (x,y)} = \frac{1}{N} \sum_{i=1}^N \text{z-score of } x_i \times \text{z-score of } y_i$$

Calculating correlations – Step 2: Calculating the correlation

$$\text{Correlation (x,y)} = \frac{1}{N} \sum_{i=1}^N \text{z-score of } x_i \times \text{z-score of } y_i$$

$$\text{Correlation (x,y)} = \frac{x_i - \bar{x}}{sd_x} \times \frac{y_i - \bar{y}}{sd_y}$$

Calculating correlations – Step 3: Interpretation

- Correlation measures *linear* association
- Correlations are between -1 and 1

Calculating correlations – Step 3: Interpretation

```
cor(log(data$GDP),data$Child.Mortality, use = "pairwise")  
  
= -0.7684907
```

Calculating correlations – Step 3: Interpretation

```
cor(data[, c("GDP", "Child.Mortality", "PolityIV")], use =  
"pairwise.complete.obs")
```

Calculating correlations – Step 3: Interpretation

```
### by hand
z_gdp <- (log(data$GDP) - mean(log(data$GDP), na.rm = T))/sd(log(data$GDP), na.rm =T)
z_CM <- (data$Child.Mortality - mean(data$Child.Mortality, na.rm = T))/sd(data$Child.Mortality, na.rm =T)
cor <- sum(z_gdp*z_CM)/(length(z_gdp)-1)
-0.7684907
```