

#### Menu

Xiao-Li Meng Department of Statistics, Harvard University

 $Big \neq Bette$ 

Motivatio

Soup

Euler Identit

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

# How Small Are Our Big Data: Turning the 2016 Surprise into a 2020 Vision

## Xiao-Li Meng Department of Statistics, Harvard University

▲ロト ▲帰 ト ▲ ヨ ト ▲ ヨ ト ● ● ● の Q ()



#### Menu

Xiao-Li Meng Department of Statistics, Harvard University

- Big ≠ Bette Motivation
- Soup
- Euler Identity
- Derivation
- Trio
- LLP
- CCES
- Assessing d.d.i
- Paradox
- Lessons

# How Small Are Our Big Data: Turning the 2016 Surprise into a 2020 Vision

Xiao-Li Meng Department of Statistics, Harvard University

- Meng and Xie (2014). I Got More Data, My Model Is More Refined, But My Estimator Is Getting Worse! Am I Just Dumb? Econometric Reviews. 33: 218-250.
- Meng (2018) Statistical Paradises and Paradoxes in Big Data (I): Law of Large Populations, Big Data Paradox, and the 2016 US Election. Annals of Applied Statistics. 2: 685-726



#### Menu

Xiao-Li Meng Department of Statistics, Harvard University

- Big ≠ Bette Motivation
- Soup
- Euler Identity
- Derivation
- Trio
- LLP
- CCES
- Assessing d.d.i
- Paradox
- Lessons

# How Small Are Our Big Data: Turning the 2016 Surprise into a 2020 Vision

Xiao-Li Meng Department of Statistics, Harvard University

- Meng and Xie (2014). I Got More Data, My Model Is More Refined, But My Estimator Is Getting Worse! Am I Just Dumb? Econometric Reviews. 33: 218-250.
- Meng (2018) Statistical Paradises and Paradoxes in Big Data (I): Law of Large Populations, Big Data Paradox, and the 2016 US Election. Annals of Applied Statistics. 2: 685-726
- Many thanks to **Stephen Ansolabehere and Shiro Kuriwaki** for the CCES (**Cooperative Congressional Election Study**) data and analysis on 2016 US election.



# "Big Data" is everywhere ...

Menu	2
Xiao-Li Mei Department	ng of
Statistics, Harvard	,
	er
	ty
CCES	
	d.i
Lessons	

▲□▶ ▲□▶ ▲ 臣▶ ▲ 臣▶ ― 臣 … のへぐ



## "Big Data" is everywhere ...

Menu

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identit Derivation Trio LLP

CCES

Assessing d.d.i

Paradox

Lessons



Photo: Sandy Zabell

◆□▶ ◆圖▶ ◆注▶ ◆注▶



# But Bigger $\Rightarrow$ Better, even for Least-Squares (LSE)





# $\mathsf{Bigger} \not\Rightarrow \mathsf{Better:} \ \mathsf{A} \ \mathsf{Mathematical} \ \mathsf{Proof}$

### Menu

Xiao-Li Meng Department of Statistics, Harvard University

## $\mathsf{Big} \neq \mathsf{Better}$

Motivation Soup Euler Identi Derivation Trio LLP CCES Assessing d.

Paradox

Lessons

## A Heteroscedastic Regression Model

$$Y_i = \beta X_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2 X_i^{\eta}), \quad i = 1, \dots, n$$



### Menu

4

### $Big \neq Better$

## A Heteroscedastic Regression Model

$$Y_i = \beta X_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2 X_i^{\eta}), \quad i = 1, \dots, n$$

## Least-squares estimator:

 $\hat{\beta}^{LSE} = \frac{\sum_{i=1}^{n} X_i Y_i}{\sum_{i=1}^{n} X_i^2}$ 



### Menu

4

Xiao-Li Meng Department of Statistics, Harvard University

### $Big \neq Better$

Motivation Soup Euler Identit Derivation Trio LLP CCES Assessing d.

### Paradox

Lessons

## A Heteroscedastic Regression Model

$$Y_i = \beta X_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2 X_i^{\eta}), \quad i = 1, \dots, n$$





#### Menu

4

Xiao-Li Meng Department of Statistics, Harvard University

### $Big \neq Better$

Motivation Soup Euler Identi Derivation Trio LLP CCES

#### Assessing d.d.i

Paradox

Lessons

## A Heteroscedastic Regression Model

$$Y_i = \beta X_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2 X_i^{\eta}), \quad i = 1, \dots, n$$



But  $\hat{\beta}^{LSE}$  is not self-efficient (Meng, 1994) when  $\eta \neq 0$ :

$$\mathbf{V}(\hat{\beta}^{LSE}|X,\theta) = \sigma^2 \frac{\sum_{i=1}^n X_i^{2+\eta}}{\left[\sum_{i=1}^n X_i^2\right]^2}$$



#### Menu

4

Xiao-Li Meng Department of Statistics, Harvard University

### $\mathsf{Big} \neq \mathsf{Better}$

Motivation Soup Euler Identit Derivation Trio LLP CCES Assessing d.

Paradox

Lessons

## A Heteroscedastic Regression Model

$$Y_i = \beta X_i + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2 X_i^{\eta}), \quad i = 1, \dots, n$$



But  $\hat{\beta}^{LSE}$  is not *self-efficient* (Meng, 1994) when  $\eta \neq 0$ :

$$\mathcal{V}(\hat{\beta}^{LSE}|X,\theta) = \sigma^2 \frac{\sum_{i=1}^{n} X_i^{2+\eta}}{\left[\sum_{i=1}^{n} X_i^2\right]^2}$$

Compare, when  $\eta = 0$ :

$$\mathcal{N}(\hat{\beta}^{LSE}|X,\theta) = \sigma^2 \frac{1}{\sum_{i=1}^n X_i^2}$$



Menu

Xiao-Li Meng Department of Statistics, Harvard University

## $\mathsf{Big} \neq \mathsf{Better}$

Motivation Soup Euler Identit Derivation Trio LLP CCES Assessing d.

Paradox

Lessons

• Those observations with large variabilities received more weight than they deserve.

▲ロト ▲帰 ト ▲ ヨ ト ▲ ヨ ト ● ● ● の Q ()



Menu

Xiao-Li Meng Department of Statistics, Harvard University

## $\mathsf{Big} \neq \mathsf{Better}$

Motivation Soup Euler Identit Derivation Trio LLP CCES

- Assessing d.d.
- Paradox

Lessons

• Those observations with large variabilities received more weight than they deserve.

Re-weight the data by  $W_i = X_i^{-\eta/2}$  (assume  $\eta$  is known)

 $W_i Y_i = \beta(W_i X_i) + \tilde{\epsilon}_i, \quad \tilde{\epsilon}_i \sim N(0, \sigma^2), \quad i = 1, \dots, n$ 



Menu

Xiao-Li Meng Department of Statistics, Harvard University

5

## $\mathsf{Big} \neq \mathsf{Better}$

Motivation Soup Euler Ident Derivation Trio LLP CCES Assessing d • Those observations with large variabilities received more weight than they deserve.

Re-weight the data by  $W_i = X_i^{-\eta/2}$  (assume  $\eta$  is known)

$$W_i Y_i = \beta(W_i X_i) + \tilde{\epsilon}_i, \quad \tilde{\epsilon}_i \sim N(0, \sigma^2), \quad i = 1, \dots, n$$

$$\hat{B}_{MLE} = \frac{\sum_{i=1}^{n} X_i^{1-\eta} Y_i}{\sum_{i=1}^{n} X_i^{2-\eta}}$$



Menu

Xiao-Li Meng Department of Statistics, Harvard University

### $Big \neq Better$

Motivation Soup Euler Ident Derivation Trio LLP CCES Assessing c • Those observations with large variabilities received more weight than they deserve.

Re-weight the data by  $W_i = X_i^{-\eta/2}$  (assume  $\eta$  is known)  $W_i Y_i = \beta(W_i X_i) + \tilde{\epsilon}_i, \quad \tilde{\epsilon}_i \sim N(0, \sigma^2), \quad i = 1, \dots, n$ 

$$\hat{P}_{MLE} = \frac{\sum_{i=1}^{n} X_{i}^{1-\eta} Y_{i}}{\sum_{i=1}^{n} X_{i}^{2-\eta}}$$

$$\mathbf{V}(\hat{\beta}_{MLE}|\mathbf{X}, \theta) = \sigma^2 \frac{1}{\sum_{i=1}^{n} X_i^{2-\eta}}$$



Menu

Xiao-Li Meng Department of Statistics, Harvard University

5

Â

## $Big \neq Better$

Motivatio Soup Euler Iden Derivatior

THO

LLP

CCES

Assessing d.d.i

Paradox

Lessons

• Those observations with large variabilities received more weight than they deserve.

Re-weight the data by  $W_i = X_i^{-\eta/2}$  (assume  $\eta$  is known)

$$W_i Y_i = \beta(W_i X_i) + \tilde{\epsilon}_i, \quad \tilde{\epsilon}_i \sim N(0, \sigma^2), \quad i = 1, \dots, n$$

$$MLE = \frac{\sum_{i=1}^{n} X_{i}^{1-\eta} Y_{i}}{\sum_{i=1}^{n} X_{i}^{2-\eta}}$$

$$\mathcal{V}(\hat{\beta}_{\mathsf{MLE}}|X,\theta) = \sigma^2 \frac{1}{\sum_{i=1}^n X_i^{2-\eta}}$$

So it is justifiable to throw away some data points if you don't know how to use them most effectively because



Menu

Xiao-Li Meng Department of Statistics, Harvard University

### $Big \neq Better$

Motivation Soup Euler Iden Derivation Trio LLP

Â

CCES

Assessing d.d.i

Paradox

Lessons

• Those observations with large variabilities received more weight than they deserve.

Re-weight the data by  $W_i = X_i^{-\eta/2}$  (assume  $\eta$  is known)

$$W_i Y_i = \beta(W_i X_i) + \tilde{\epsilon}_i, \quad \tilde{\epsilon}_i \sim N(0, \sigma^2), \quad i = 1, \dots, n$$

$$MLE = \frac{\sum_{i=1}^{n} X_{i}^{1-\eta} Y_{i}}{\sum_{i=1}^{n} X_{i}^{2-\eta}}$$

$$\mathcal{V}(\hat{\beta}_{\mathsf{MLE}}|X,\theta) = \sigma^2 \frac{1}{\sum_{i=1}^n X_i^{2-\eta}}$$

So it is justifiable to throw away some data points if you don't know how to use them most effectively because

When the optimal  $W_i$ 's have large variation, setting small  $W_i$ 's to zero better approximates the optimal weighting scheme than "blindly" using equal weights.



## Don't sue me ...

#### Menu

Xiao-Li Meng Department of Statistics, Harvard University

6

### $\mathsf{Big} \neq \mathsf{Better}$

- Motivatio Soup Fuler Ide
- Derivation
- Trio
- LLP
- CCES
- Assessing d.d.i
- Paradox
- Lessons

## WHEN IT SEEMS DESIRABLE TO IGNORE DATA

#### Herman Chernoff

Massachusetts Institute of Technology

#### ABSTRACT

An experiment designed to detect the relative motion of two astronomical objects raised the problem of testing, against shift alternatives, the hypothesis H<sub>0</sub> that two energy distributions are equivalent. The relevant data consist of independent Poisson counts  $X_{i,j}$  with means  $\lambda_j p_{i,j} T_{i,j}$  where  $\lambda_j$  is the intensity of radiation from the jth object,  $p_{i,j}$  is the probability that a random photon from the jth object has energy in a small interval centered about  $e_i$ , and  $T_{i,j}$  is the time duration allocated to the count  $X_{i,j}$ . The hypothesis H<sub>0</sub> implies that  $p_{i,1} = p_{i,2}$  for  $i = 1, 2, \dots, m$ .

A natural test uses the statistic  $E_i(\hat{p}_{12} - \hat{p}_{11})$  where the  $\hat{p}_{1j}$  are estimates of  $p_{1j}$ . For intervals where the  $p_{1j}$  were anticipated to be small, the experimenter choice small  $T_{ij}$ , values



# From Robins and Wang (2000, Biometrika)





# Data Quality vs Quantity: Motivating questions

#### Menu

8

Xiao-Li Meng Department of Statistics, Harvard University

#### $Big \neq Better$

## Motivation

- Soup
- Destantes
- Trio
- LLP
- CCES
- Assessing d.d.i
- Paradox
- Lessons

• We know that a 5% random sample is better than a 5% non-random sample in measurable ways (e.g., bias, predictive power).



# Data Quality vs Quantity: Motivating questions

#### Menu

8

- Xiao-Li Meng Department of Statistics, Harvard University
- Big ≠ Bettei
- Motivation
- Soup
- Euler Identity
- Derivation
- Trio
- LLP
- CCES
- Assessing d.d.i
- Paradox
- Lessons

- We know that a 5% random sample is better than a 5% non-random sample in measurable ways (e.g., bias, predictive power).
- But is an 80% non-random sample "better" than a 5% random sample in measurable terms? 90%?
  95%? 99%? (Wu, 2012, Seminar at Harvard Statistics)

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ つ へ つ ・



# Data Quality vs Quantity: Motivating questions

#### Menu

- Xiao-Li Meng Department of Statistics, Harvard University
- Big ≠ Bettei
- Motivation
- Soup
- Euler Identity
- Derivation
- Trio
- LLP
- CCES
- Assessing d.d.i
- Paradox
- Lessons

- We know that a 5% random sample is better than a 5% non-random sample in measurable ways (e.g., bias, predictive power).
- But is an 80% non-random sample "better" than a 5% random sample in measurable terms? 90%?
  95%? 99%? (Wu, 2012, Seminar at Harvard Statistics)
- "Which one should we trust more: a 1% survey with 60% response rate or a non-probabilistic dataset covering 80% of the population?" (Keiding and Louis, 2015, Joint Statistical Meetings; and *JRSSB*, 2016)



イロト 不得下 不良下 不良下

3

990



Xiao-Li Meng Department of Statistics, Harvard University

q

Big ≠ Better

Motivation Soup

Euler Identit

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

## • Law of Large Numbers: Jakob Bernoulli (1713)



◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○

Menu

q

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Bette

Motivation

Soup

Tele

\_ \_ \_ \_

Assessing d.d.i

Paradox

Lessons

• Law of Large Numbers: Jakob Bernoulli (1713)

 Central Limit Theorem: Abraham de Moivre (1733): error ∝ 1/√n : n − sample size



Menu

Xiao-Li Meng Department of Statistics, Harvard University

q

Big ≠ Bette

## Motivation

Soup

Euler Identity

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

- Law of Large Numbers: Jakob Bernoulli (1713)
- Central Limit Theorem: Abraham de Moivre (1733): error ∝ 1/√n : n − sample size

## • Survey Sampling:

• Graunt (1662); Laplace (1882)



Menu

Xiao-Li Meng Department of Statistics, Harvard University

q

Big ≠ Bette

## Motivation

Soup

Euler Identity

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

- Law of Large Numbers: Jakob Bernoulli (1713)
- Central Limit Theorem: Abraham de Moivre (1733): error ∝ 1/√n : n − sample size

## • Survey Sampling:

- Graunt (1662); Laplace (1882)
- The "intellectually violent revolution" in 1895 by Anders Kiær, Statistics Norway



Menu

q

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Bette

## Motivation

Soup

Euler Identity

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

- Law of Large Numbers: Jakob Bernoulli (1713)
- Central Limit Theorem: Abraham de Moivre (1733): error ∝ 1/√n : n − sample size

## • Survey Sampling:

- Graunt (1662); Laplace (1882)
- The "intellectually violent revolution" in 1895 by Anders Kiær, Statistics Norway





Menu

q

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Bette

## Motivation

Soup

Euler Identity

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

- Law of Large Numbers: Jakob Bernoulli (1713)
- Central Limit Theorem: Abraham de Moivre (1733): error ∝ 1/√n : n − sample size

## • Survey Sampling:

- Graunt (1662); Laplace (1882)
- The "intellectually violent revolution" in 1895 by Anders Kiær, Statistics Norway





Menu

Xiao-Li Meng Department of Statistics, Harvard University

q

 $\mathsf{Big} \neq \mathsf{Bette}$ 

## Motivation

Soup

Euler Identity

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

- Law of Large Numbers: Jakob Bernoulli (1713)
- Central Limit Theorem: Abraham de Moivre (1733): error ∝ 1/√n : n − sample size

## • Survey Sampling:

- Graunt (1662); Laplace (1882)
- The "intellectually violent revolution" in 1895 by Anders Kiær, Statistics Norway



• Landmark paper: Jerzy Neyman (1934)



### Menu

Xiao-Li Meng Department of Statistics, Harvard University

q

Big ≠ Bette

## Motivation

Soup

Euler Identity

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

- Law of Large Numbers: Jakob Bernoulli (1713)
- Central Limit Theorem: Abraham de Moivre (1733): error ∝ <sup>1</sup>/<sub>√n</sub> : n − sample size

## • Survey Sampling:

- Graunt (1662); Laplace (1882)
- The "intellectually violent revolution" in 1895 by Anders Kiær, Statistics Norway



- Landmark paper: Jerzy Neyman (1934)
- The "revolution" lasted about 50 years (Jelke Bethlehem, 2009)



### Menu

Xiao-Li Meng Department of Statistics, Harvard University

q

 $\mathsf{Big} \neq \mathsf{Bette}$ 

## Motivation

Soup

Euler Identity

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

- Law of Large Numbers: Jakob Bernoulli (1713)
- Central Limit Theorem: Abraham de Moivre (1733): error ∝ <sup>1</sup>/<sub>√n</sub> : n − sample size

## • Survey Sampling:

- Graunt (1662); Laplace (1882)
- The "intellectually violent revolution" in 1895 by Anders Kiær, Statistics Norway



- Landmark paper: Jerzy Neyman (1934)
- The "revolution" lasted about 50 years (Jelke Bethlehem, 2009)
- First implementation in US Census: 1940 led by Morris Hansen





#### Menu 10

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Bette

Motivation

## Soup

Euler Identit

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

• Think about tasting soup ...



#### Menu 10

Xiao-Li Meng Department of Statistics, Harvard University

 $Big \neq Bette$ 

Motivation

## Soup

Euler Identit

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

- Think about tasting soup ...
- Stir it well, then a few bits are sufficient regardless of the size of the container!





#### Menu :

- Xiao-Li Meng Department of Statistics, Harvard University
- $Big \neq Bette$
- Motivation

## Soup

- Euler Identit
- Derivation
- Trio
- LLP
- CCES
- Assessing d.d.i
- Paradox
- Lessons

- Think about tasting soup ...
- Stir it well, then a few bits are sufficient regardless of the size of the container!





### Menu 10

- Xiao-Li Meng Department of Statistics, Harvard University
- $\mathsf{Big} \neq \mathsf{Bette}$
- Motivation

## Soup

- Euler Identity Derivation Trio LLP CCES
- Assessing d.d.i Paradox
- Lessons

- Think about tasting soup ...
- Stir it well, then a few bits are sufficient regardless of the size of the container!



A D F A B F A B F A B F





### Menu 10

Xiao-Li Meng Department of Statistics, Harvard University

 $\mathsf{Big} \neq \mathsf{Bette}$ 

Motivation

## Soup

Euler Identity Derivation Trio LLP CCES

Assessing d.d.i Paradox

I diadox

Lessons

- Think about tasting soup ...
- Stir it well, then a few bits are sufficient regardless of the size of the container!


















CREATIVEDAFFODIL.ETSY.COM



・ロト ・ 同ト ・ ヨト ・ ヨト

3



### Menu

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identity

Derivatio

Irio

LLP

CCES

Assessing d.d.i

Paradox

Lessons



CREATIVEDAFFODIL ETSY.COM



▲ロト ▲帰 ト ▲ ヨ ト ▲ ヨ ト ● ● ● の Q ()

• 5 most fundamental numbers in mathematics:

 $0, 1, e, \pi, i = \sqrt{-1}$ 





Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identity

Derivatio

Irio

LLP

CCES

Assessing d.d.

Paradox

Lessons



CREATIVEDAFFODIL ETSY.COM



• 5 most fundamental numbers in mathematics:

D, 1, 
$$e, \pi, i$$
 =  $\sqrt{-1}$ 

• The unexpected one:  $i = \sqrt{-1}$ 



Menu 12	What are the five most fundamental symbols in Statistics?	
Xiao-Li Meng Department of Statistics, Harvard University	• $\mu$ : Average/Mean Ave $\{X_j, j = 1,\}$	
$Big \neq Better$		
Motivation		
Soup		
Euler Identity		
Derivation		
Trio		
LLP		
CCES		
Assessing d.d.i		
Paradox		
Lessons		

◆ロト ◆昼 ト ◆臣 ト ◆臣 ト ○日 ○ のへで



# What are the five most fundamental symbols in Statistics? Menu Average/Mean Ave $\{X_i, j = 1, ...\}$ • μ: Department of $\sqrt{\operatorname{Ave}\{(X_i - \mu)^2\}}$ Standard Deviation • $\sigma$ : Euler Identity



### Menu

#### Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup

### Euler Identity Derivation Trio LLP

CCES

Assessing d.d.

Paradox

Lessons

## What are the five most fundamental symbols in Statistics?

- $\mu$ : Average/Mean
- $\sigma$ : Standard Deviation
- $\rho$ : Correlation

 $\sqrt{Ave\{(X_i - \mu)^2\}}$  $\operatorname{Ave}\left(\frac{X_j}{\sigma_x}\frac{Y_j}{\sigma_y}\right) - \operatorname{Ave}\left(\frac{X_j}{\sigma_x}\right)\operatorname{Ave}\left(\frac{Y_j}{\sigma_y}\right)$ 

Ave $\{X_i, j = 1, ...\}$ 



### Menu

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup

Euler Identity Derivation Trio LLP

CCES

Assessing d.d.i

Paradox

Lessons

## What are the five most fundamental symbols in Statistics?

Av

- $\mu$ : Average/Mean
- $\sigma$ : Standard Deviation
- $\rho$ : Correlation

$$\sqrt{\operatorname{Ave}\{(X_j - \mu)^2\}} \\ = \left(\frac{X_j}{\sigma_x}\frac{Y_j}{\sigma_y}\right) - \operatorname{Ave}(\frac{X_j}{\sigma_x})\operatorname{Ave}(\frac{Y_j}{\sigma_y})$$

Ave $\{X_i, j = 1, ...\}$ 

 $\left[ A \left( \left( X \right) \right)^{2} \right]$ 

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○

• n: Sample Size



### Menu

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation

Soup

Euler Identity
Derivation

LLP

CCES

Assessing d.d.i

Paradox

Lessons

## What are the five most fundamental symbols in Statistics?

- $\mu$ : Average/Mean Ave $\{X_j, j = 1, ...\}$
- $\sigma$ : Standard Deviation

Population Size

•  $\rho$ : Correlation

• N:

The unexpected one ...

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○

 $\operatorname{Ave}\left(\frac{X_{j}}{\sigma_{x}}\frac{Y_{j}}{\sigma_{y}}\right) - \operatorname{Ave}\left(\frac{X_{j}}{\sigma_{x}}\right)\operatorname{Ave}\left(\frac{Y_{j}}{\sigma_{y}}\right)$ 

 $\sqrt{Ave\{(X_i - \mu)^2\}}$ 



### Menu

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation

Soup

Euler Identity
Derivation

LLP

CCES

Assessing d.d.i

Paradox

Lessons

## What are the five most fundamental symbols in Statistics?

- $\mu$ : Average/Mean Ave $\{X_j, j = 1, ...\}$
- $\sigma$ : Standard Deviation

Population Size

•  $\rho$ : Correlation

• N:

The unexpected one ...

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○

 $\operatorname{Ave}\left(\frac{X_{j}}{\sigma_{x}}\frac{Y_{j}}{\sigma_{y}}\right) - \operatorname{Ave}\left(\frac{X_{j}}{\sigma_{x}}\right)\operatorname{Ave}\left(\frac{Y_{j}}{\sigma_{y}}\right)$ 

 $\sqrt{Ave\{(X_i - \mu)^2\}}$ 



#### Menu

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation

Soup

Euler Identity

Trio

LLP

CCES

Assessing d.d.

Paradox

Lessons

## What are the five most fundamental symbols in Statistics?

Ave $\{X_i, j = 1, ...\}$ 

 $\sqrt{Ave\{(X_i - \mu)^2\}}$ 

 $\operatorname{Ave}\left(\frac{X_j}{\sigma_x}\frac{Y_j}{\sigma_y}\right) - \operatorname{Ave}\left(\frac{X_j}{\sigma_x}\right)\operatorname{Ave}\left(\frac{Y_j}{\sigma_y}\right)$ 

- $\mu$ : Average/Mean
- $\sigma$ : Standard Deviation
- $\rho$ : Correlation
- n: Sample Size
- N: Population Size The unexpected one ...

### The Most Beautiful Statistical Identity?

$$\hat{\mu}_n - \mu_N = \hat{\rho}\sigma\sqrt{\frac{N-n}{n}}$$

・ロト ・西ト ・田ト ・田ト ・日・ シタマ



Menu 13
---------

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Bette Motivation Soup Euler Identit

### Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

• *n*: number of respondents to an election survey

▲ロト ▲帰 ト ▲ ヨ ト ▲ ヨ ト ● ● ● の Q ()



Menu

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Bette Motivation Soup Fuler Identit

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

• *n*: number of respondents to an election survey

▲ロト ▲帰 ト ▲ ヨ ト ▲ ヨ ト ● ● ● の Q ()

• N: number of (actual) voters in US



Menu 13

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identit<u>y</u>

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

- *n*: number of respondents to an election survey
- N: number of (actual) voters in US
- $X_j = 1$ : plan to vote for Trump;  $X_j = 0$  otherwise



#### Menu 13

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup

#### Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

- *n*: number of respondents to an election survey
- N: number of (actual) voters in US
- $X_j = 1$ : plan to vote for Trump;  $X_j = 0$  otherwise
- $R_j = 1$ : report (honestly) voting plan;  $R_j = 0$  otherwise



#### Menu 13

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup

#### Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

- *n*: number of respondents to an election survey
- N: number of (actual) voters in US
- $X_j = 1$ : plan to vote for Trump;  $X_j = 0$  otherwise
- $R_j = 1$ : report (honestly) voting plan;  $R_j = 0$  otherwise



Menu 13

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identity

Derivation

Trio LLP CCES Assessing Paradox

Lessons

- *n*: number of respondents to an election survey
- N: number of (actual) voters in US
- $X_j = 1$ : plan to vote for Trump;  $X_j = 0$  otherwise
- $R_j = 1$ : report (honestly) voting plan;  $R_j = 0$  otherwise

# Estimatinng Trump's share: $\mu_{\scriptscriptstyle N}={\sf Ave}(X_j)$ by sample average:

$$\hat{\mu}_n = \frac{R_1 X_1 + \ldots + R_N X_N}{n} = \frac{\operatorname{Ave}(R_j X_j)}{\operatorname{Ave}(R_j)}$$



Menu

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identity

### Derivation

Trio LLP

Assessing d.d

Paradox

Lessons

- *n*: number of respondents to an election survey
- N: number of (actual) voters in US
- $X_j = 1$ : plan to vote for Trump;  $X_j = 0$  otherwise
- $R_j = 1$ : report (honestly) voting plan;  $R_j = 0$  otherwise

## Estimationg Trump's share: $\mu_N = Ave(X_j)$ by sample average:

$$\hat{\mu}_n = \frac{R_1 X_1 + \ldots + R_N X_N}{n} = \frac{\operatorname{Ave}(R_j X_j)}{\operatorname{Ave}(R_j)}$$

### Actual estimation error

$$\hat{\mu}_{n} - \mu_{N} = \frac{\operatorname{Ave}(R_{j}X_{j})}{\operatorname{Ave}(R_{j})} - \operatorname{Ave}(X_{j})$$
$$= \left[\frac{\operatorname{Ave}(R_{j}X_{j}) - \operatorname{Ave}(R_{j})\operatorname{Ave}(X_{j})}{\sigma_{R}\sigma_{X}}\right] \times \frac{\sigma_{R}}{\operatorname{Ave}(R_{j})} \times \sigma_{X}$$



## Data quality, quantity, and uncertainty

#### Menu 14

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identity Derivation **Trio** 

LLP

CCES

Assessing d.d.i

Paradox

Lessons

## Because $\sigma_R^2 = f(1-f)$ , $f = Ave\{R_j\} = \frac{n}{N}$ , we have

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○

$$\text{Error} = \underbrace{\hat{\rho}_{R,X}}_{\text{Data Quality}} \times$$



## Data quality, quantity, and uncertainty

### Menu 14 Xiao-Li Meng Department of Statistics, Harvard University Big $\neq$ Better Motivation Soup

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○

Trio

CCES

Assessing d.d.i

Paradox

Lessons



Trio

## Data quality, quantity, and uncertainty



◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○



#### Menu 15

Xiao-Li Meng Department of Statistics, Harvard University

- $\operatorname{Big} \neq \operatorname{Bette}$
- Motivatio
- Soup
- Euler Identity
- Derivation

### Trio

- LLP
- CCES
- Assessing d.d.i
- Paradox
- Lessons

### Mean Squared Error (MSE)

$$\operatorname{MSE}(\hat{\mu}_n) = \mathsf{E}_R(\hat{\rho}^2) \times \frac{N-n}{n} \times \sigma_x^2$$



### Menu 15

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identity Derivation

### Trio

LLP

CCES

Assessing d.d

Paradox

Lessons

# Mean Squared Error (MSE)

$$\operatorname{MSE}(\hat{\mu}_n) = \mathsf{E}_R(\hat{\rho}^2) \times \frac{N-n}{n} \times \sigma_x^2$$

## Data Defect Index (d.d.i): $D_I = E_R(\hat{\rho}^2)$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□▶



### Menu 15

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identity Derivation

### Trio

LLP

- CCES
- Assessing d.d.i
- Paradox
- Lessons

### Mean Squared Error (MSE)

$$\operatorname{MSE}(\hat{\mu}_n) = \mathsf{E}_R(\hat{\rho}^2) \times \frac{N-n}{n} \times \sigma_x^2$$

## Data Defect Index (d.d.i): $D_I = \mathsf{E}_{\mathcal{R}}(\hat{\rho}^2)$

• For Simple Random Sample (SRS):  $D_I = (N-1)^{-1}$ 

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト … ヨ



### Menu 15

Xiao-Li Meng Department of Statistics, Harvard University

- Big ≠ Bette Motivation
- Soup
- Euler Identity
- Derivation

### Trio

- LLP
- CCES
- Assessing d.d.i
- Paradox
- Lessons

### Mean Squared Error (MSE)

$$\operatorname{MSE}(\hat{\mu}_n) = \mathsf{E}_R(\hat{\rho}^2) \times \frac{N-n}{n} \times \sigma_x^2$$

## Data Defect Index (d.d.i): $D_I = \mathsf{E}_{\mathcal{R}}(\hat{\rho}^2)$

• For Simple Random Sample (SRS):  $D_I = (N-1)^{-1}$ 

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト … ヨ

• For probabilistic samples in general:  $D_I \propto N^{-1}$ 



#### Menu 15

Xiao-Li Meng Department of Statistics, Harvard University

- Big ≠ Bette Motivation
- Soup
- Euler Identity
- Derivation
- Trio
- LLP
- CCES
- Assessing d.d.i
- Paradox
- Lessons

### Mean Squared Error (MSE)

$$\operatorname{MSE}(\hat{\mu}_n) = \mathsf{E}_R(\hat{\rho}^2) \times \frac{N-n}{n} \times \sigma_x^2$$

## Data Defect Index (d.d.i): $D_I = E_R(\hat{\rho}^2)$

- For Simple Random Sample (SRS):  $D_I = (N-1)^{-1}$
- For probabilistic samples in general:  $D_I \propto N^{-1}$
- Deep trouble when  $D_I$  does not vanish with  $N^{-1}$ ;
- or equivalently when  $\hat{
  ho}$  does not vanish with  $N^{-1/2}$  ...

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○



# A Law of Large Populations (LLP)

### Menu

16

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Bette Motivation Soup Euler Identi Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

## If $\rho = \mathsf{E}_{R}(\hat{\rho}) \neq 0$ , then on average, the relative error $\uparrow \sqrt{N}$ :

$$"Z-\text{score"} \equiv \frac{\text{Actual Error}}{\text{Benchmark SRS Standard Error}} = \sqrt{N-1}\hat{\rho}$$



# A Law of Large Populations (LLP)

### Menu

Xiao-Li Meng Department of Statistics, Harvard University

16

Big ≠ Bette Motivation Soup Euler Identi Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

## If $\rho = \mathsf{E}_{R}(\hat{\rho}) \neq 0$ , then on average, the relative error $\uparrow \sqrt{N}$ :

"Z-score" 
$$\equiv \frac{\text{Actual Error}}{\text{Benchmark SRS Standard Error}} = \sqrt{N-1}\hat{\rho}$$

## The (lack-of) design effect (Deff)

$$\mathrm{Deff} = \frac{\mathrm{MSE}}{\mathrm{Benchmark \ SRS \ MSE}} = (N-1)D_I$$

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○



# A Law of Large Populations (LLP)

### Menu 16

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identity Derivation Trio

ПР

CCES

Assessing d.d.i

Paradox

Lessons

## If $\rho = \mathsf{E}_{R}(\hat{\rho}) \neq 0$ , then on average, the relative error $\uparrow \sqrt{N}$ :

"Z-score" 
$$\equiv \frac{\text{Actual Error}}{\text{Benchmark SRS Standard Error}} = \sqrt{N-1}\hat{\rho}$$

### The (lack-of) design effect (Deff)

$$\mathrm{Deff} = rac{\mathrm{MSE}}{\mathrm{Benchmark \ SRS \ MSE}} = (N-1)D_I$$

The *Effective Sample Size*  $n_{\rm eff}$  of a "Big Data" set

Equate its MSE to that from a SRS with size  $n_{\rm eff}$ :

$$D_{I}\left[\frac{1-f}{f}\right]\sigma^{2} = \frac{1}{N-1}\left[\frac{N-n_{\text{eff}}}{n_{\text{eff}}}\right]\sigma^{2}$$

◆□> ◆□> ◆豆> ◆豆> ・豆 ・ のへで



# Gaining 2020 Vision: Assessing the behavioral $\hat{\rho}$

#### Menu

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identity Derivation Trio LLP CCES Assessing d.d.i

Paradox

Lessons

### CCES: Cooperative Congressional Election Study

(Conducted by Stephen Ansolabehere, Brian Schaffner, Sam Luks, Douglas Rivers on **Oct 4** - **Nov 6**, **2016** (YouGov); Analysis assisted by Shiro Kuriwaki)



Raw Sample: 64,600 Voting Adj: 48,106 Validated: 34,156

### Reasonable predictions for Clinton's Vote Share



# Gross under-prediction/reporting of Trump's Share

#### Menu 18

## CCES: Cooperative Congressional Election Study



Raw Sample: 64,600 Voting Adj: 48,106 Validated: 34,156

There are many "undecided" ...



#### Menu 19

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identity Derivation Trio LLP

Assessing d.d.i

### Let $\mu_{N}$ be the true share, and $\hat{\mu}_{n}$ the estimated share. Then

$$\hat{
ho} = rac{\hat{\mu}_n - \mu_N}{\sqrt{rac{N-n}{n}\sigma^2}}, \quad \& \quad \sigma^2 = \mu_N(1-\mu_N)$$



#### Menu 19

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identity Derivation Trio LLP

### Assessing d.d.i

Paradox

Lessons

### Let $\mu_N$ be the true share, and $\hat{\mu}_n$ the estimated share. Then

$$\hat{\rho} = \frac{\hat{\mu}_n - \mu_N}{\sqrt{\frac{N-n}{n}\sigma^2}}, \quad \& \quad \sigma^2 = \mu_N (1 - \mu_N)$$





#### Menu 19

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identit Derivation Trio LLP

CCES

Assessing d.d.i

Paradox

Lessons

### Let $\mu_N$ be the true share, and $\hat{\mu}_n$ the estimated share. Then

$$\hat{
ho} = rac{\hat{\mu}_n - \mu_N}{\sqrt{rac{N-n}{n}\sigma^2}}, \quad \& \quad \sigma^2 = \mu_N (1 - \mu_N)$$





#### Menu 19

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Better Motivation Soup Euler Identit Derivation Trio LLP

CCES

Assessing d.d.i

Paradox

Lessons

### Let $\mu_N$ be the true share, and $\hat{\mu}_n$ the estimated share. Then

$$\hat{
ho} = rac{\hat{\mu}_n - \mu_N}{\sqrt{rac{N-n}{n}\sigma^2}}, \quad \& \quad \sigma^2 = \mu_N(1-\mu_N)$$



• Problem: Voter validation is done through matching algorithms and it is not fool-proof, and it may introduce additional *selection bias*.



## What's the implication of $\hat{\rho} = -0.005$ ?

Menu 20		
Xiao-Li Meng Department of Statistics, Harvard University		
$Big \neq Better$		
Motivation		
Soup		
Euler Identity		
Derivation		
Trio		
LLP		
CCES		
Assessing d.d.i		
Paradox		
Lessons		

• Many (major) election survey results were published daily for several months before Nov 8, 2016;


Menu 20

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Bette Motivation Soup Euler Identi1

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

- Many (major) election survey results were published daily for several months before Nov 8, 2016;
- Roughly amounts to having opinions from (up to) f = 1% of US voting eligible population: n ≈ 2, 300, 000;
- Equivalent to about 2,300 surveys of 1,000 respondents each.



Menu 20

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Bette Motivation

Soup

Euler Identit

Derivation

Trio

ПР

CCES

Assessing d.d.i

Paradox

Lessons

- Many (major) election survey results were published daily for several months before Nov 8, 2016;
- Roughly amounts to having opinions from (up to) f = 1% of US voting eligible population: n ≈ 2,300,000;
- Equivalent to about 2,300 surveys of 1,000 respondents each.

When 
$$\hat{
ho} = -0.005 = -1/200, D_I = 1/40000$$
, and hence $n_{
m eff} = rac{f}{1-f}rac{1}{D_I} = rac{1}{99} imes 40000 pprox 404!$ 



Menu 20

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Bette Motivation

Soup

Euler Identity

Derivation

Trio

ПР

CCES

Assessing d.d.i

Paradox

Lessons

- Many (major) election survey results were published daily for several months before Nov 8, 2016;
- Roughly amounts to having opinions from (up to) f = 1% of US voting eligible population: n ≈ 2,300,000;
- Equivalent to about 2,300 surveys of 1,000 respondents each.

When 
$$\hat{\rho} = -0.005 = -1/200$$
,  $D_I = 1/40000$ , and hence  
 $n_{\rm eff} = rac{f}{1-f}rac{1}{D_I} = rac{1}{99} imes 40000 pprox 404!$ 

• A 99.98% reduction in *n*, caused by  $\hat{\rho} = -0.005$ .



Menu 20

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Bette Motivation Soup

Euler Identit

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

- Many (major) election survey results were published daily for several months before Nov 8, 2016;
- Roughly amounts to having opinions from (up to) f = 1% of US voting eligible population: n ≈ 2,300,000;
- Equivalent to about 2,300 surveys of 1,000 respondents each.

When 
$$\hat{
ho} = -0.005 = -1/200, D_I = 1/40000$$
, and hence $n_{
m eff} = rac{f}{1-f}rac{1}{D_I} = rac{1}{99} imes 40000 pprox 404!$ 

- A 99.98% reduction in *n*, caused by  $\hat{\rho} = -0.005$ .
- Butterfly Effect due to Law of Large Populations (LLP)

Relative Error =  $\sqrt{N-1}\hat{\rho}$ 



# Visulizing LLP: Actual Coverage for Clinton





イロト イポト イヨト イヨト



# Visulizing LLP: Actual Coverage for Trump



◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 = つへで



# The Big Data Paradox:

Menu 23

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Bett

Motivatic

Soup

Euler Identity

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

If we do not pay attention to data quality, then

# The bigger the data, the surer we fool ourselves.

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○



#### Menu 24

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Bette

Motivatio

Soup

Euler Identit

Derivation

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

• Lesson 1: What matters most is the quality, not the quantity.



#### Menu 24

Xiao-Li Meng Department of Statistics, Harvard University

Big ≠ Bette Motivation Soup Euler Identi

Trio

LLP

CCES

Assessing d.d.i

Paradox

Lessons

- Lesson 1: What matters most is the quality, not the quantity.
- Lesson 2: Don't ignore seemingly tiny probabilistic datasets when combining data sources.

▲ロト ▲帰 ト ▲ ヨ ト ▲ ヨ ト ● ● ● の Q ()



### Menu 24

Xiao-Li Meng Department of Statistics, Harvard University

- Big ≠ Bette
- Motivatior
- Soup
- Euler Identit
- Derivation
- Trio
- LLP
- CCES
- Assessing d.d.i
- Paradox
- Lessons

- Lesson 1: What matters most is the quality, not the quantity.
- Lesson 2: Don't ignore seemingly tiny probabilistic datasets when combining data sources.
- Lesson 3: Watch the relative size, not the absolute size.

▲ロト ▲帰 ト ▲ ヨ ト ▲ ヨ ト ● ● ● の Q ()



#### Menu 24

Xiao-Li Meng Department of Statistics, Harvard University

- Big ≠ Bette
- Motivatior
- Soup
- Euler Identit
- Derivation
- Trio
- пр
- CCES
- Assessing d.d.i
- Paradox
- Lessons

- Lesson 1: What matters most is the quality, not the quantity.
- Lesson 2: Don't ignore seemingly tiny probabilistic datasets when combining data sources.
- Lesson 3: Watch the relative size, not the absolute size.
- Lesson 4: Probabilistic sampling is an extremely powerful tool to ensure data quality, but it is not the only strategy.

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○



#### Menu 24

Xiao-Li Meng Department of Statistics, Harvard University

- Big ≠ Bette
- Motivatior
- Soup
- Euler Identit
- Derivation
- Trio
- LLP
- CCES
- Assessing d.d.
- Paradox
- Lessons

- Lesson 1: What matters most is the quality, not the quantity.
- Lesson 2: Don't ignore seemingly tiny probabilistic datasets when combining data sources.
- Lesson 3: Watch the relative size, not the absolute size.
- Lesson 4: Probabilistic sampling is an extremely powerful tool to ensure data quality, but it is not the only strategy.

• Lesson 5: We may all have had too much "confidence" in big size ...



### ... and learning from real experts ...

#### Menu 25

Department of

Lessons

### 19 things we learned from the 2016 election\*

Andrew Gelman<sup>†</sup> Julia Azari<sup>‡</sup> 12 July 2017

We can all agree that the presidential election result was a shocker. According to news reports, even the Trump campaign team was stunned to come up a winner.

So now seems like a good time to go over various theories floating around in political science and political reporting and see where they stand, now that this turbulent political year has drawn to a close. In the present article, we go through several things that we as political observers and political scientists have learned from the election, and then discuss implications for the future.

### The shock

Immediately following the election there was much talk about the failure of the polls: Hillary Clinton was seen as the clear favorite for several months straight, and then she lost. After all the votes were counted, though, the view is slightly different; by election eve, the national polls were giving Clinton 52 or 53% of the two-party vote, and she ended up receiving 51%. An error of 2 percentage points is no great embarrassment.

The errors in the polls were, however, not uniform. As Figures 1 and 2 show, the Republican candidate outperformed by about 5% in highly Republican states, 2% in swing states, and not at all, on average, in highly Democratic states. This was unexpected in part because, in other recent elections, the errors in poll-based forecasts did not have this sort of structure. In 2016, though, Donald Trump won from his better-than-expected performance in Wisconsin, Michigan, North Carolina, Pennsylvania, and several other swing states.

Trump's win in the general election, and the corresponding success of Republican candidates for the U.S. Senate, then raises two questions: (1) What did the polls get wrong in these key states?, (2) How did Trump and his fellow Republicans do so well? The first is a question about survey respondents, the second a question about voters.

Going backward in time from the election-day shocker, there is the question of how Trump, as a =