

Menu 1 Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

CCES

Assessing d.d.i

Paradox

Lessons

Statistical Paradises and Paradoxes in Big Data (I): Law of Large Populations, Big Data Paradox, and the 2016 US Election

Xiao-Li Meng Department of Statistics, Harvard University



Menu

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

CCES

Assessing d.d.i

Paradox

Lessons

Statistical Paradises and Paradoxes in Big Data (I): Law of Large Populations, Big Data Paradox, and the 2016 US Election

Xiao-Li Meng Department of Statistics, Harvard University

 Meng (2018). Annals of Applied Statistics, No. 2, 685-726. https://statistics.fas.harvard.edu/people/xiao-li-meng



Menu

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

" Trio" Identity

Trio

LLP

What's Big

CCES

Assessing d.d.i

Paradox

Lessons

Statistical Paradises and Paradoxes in Big Data (I): Law of Large Populations, Big Data Paradox, and the 2016 US Election

Xiao-Li Meng Department of Statistics, Harvard University

- Meng (2018). Annals of Applied Statistics, No. 2, 685-726. https://statistics.fas.harvard.edu/people/xiao-li-meng
- Many thanks to Stephen Ansolabehere and Shiro Kuriwaki for the CCES (Cooperative Congressional Election Study) data and analysis on 2016 US election.



Motivating questions

Menu

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup " Trio

Trio

LLP

What's Big

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).



Motivating questions

Menu

2

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

- Soup
- "Trio" Identity
- Trio
- LLP
- What's Big
- 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%? (Jeremy Wu of US Census Bureau, 2012, Seminar at Harvard Statistics)

▲ロ ▶ ▲ □ ▶ ▲ □ ▶ ▲ □ ▶ ▲ □ ▶ ● ○ ○ ○



Motivating questions

Menu

2

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

- Soup
- "Trio" Identity
- Trio
- LLP
- What's Big
- 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%? (Jeremy Wu of US Census Bureau, 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)



・ロト ・ 同ト ・ ヨト ・ ヨト

= 900



Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup " Trio"

Trio

LLP

What's Big

CCES

Assessing d.d.i

Paradox

Lessons

• Law of Large Numbers: Jakob Bernoulli (1713)



◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○

Menu

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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



Menu

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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

3

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Bi

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

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Bi

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

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Bi

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

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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

3

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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)
- The "revolution" lasted about 50 years (Jelke Bethlehem, 2009)



Menu

3

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

" Trio" Identity

Trio

LLP

What's Big

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





Menu

Xiao-Li Meng Department of Statistics, Harvard University

4

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

CCES

Assessing d.d.i

Paradox

Lessons

• Think about tasting soup ...



Menu

4

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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

4

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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

4

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity Trio LLP What's Big? CCES Assessing d.d. Paradox

Lessons

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







Menu

4

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity Trio LLP What's Big? CCES Assessing d.d. Paradox

Lessons

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









	-	
- N -	lon	

5

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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

Motivatior

Soup

"Trio" Identity

Trio

LLP

What's Big

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

Xiao-Li Meng Department of Statistics, Harvard University

5

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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

5

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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

5

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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

Xiao-Li Meng Department of Statistics, Harvard University

5

Motivatior

Soup

" Trio" Identity Trio

LLP

What's Big

CCES

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

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

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



Menu

Xiao-Li Meng Department of Statistics, Harvard University

5

Motivation

Soup

"Trio" Identity

LLP

What's Big

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

Estimatinng 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

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



Data quality, quantity, and uncertainty

Menu

Xiao-Li Meng Department of Statistics, Harvard University

6

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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 6 Xiao-Li Meng Department of Statistics, Harvard University Motivation Soup Trio" Identity Menu 6 Because $\sigma_R^2 = f(1 - f), f = Ave\{R_j\} = \frac{n}{N}$, we have $\sum_{\substack{n \in \mathbb{N}, X \\ Data Quality}} \times \sqrt{\frac{N - n}{n}} \times \sqrt{\frac{n}{n}}$

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○

Trio

LLP

What's Big

CCES

Assessing d.d.

Paradox

Lessons



Data quality, quantity, and uncertainty



◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○

What's B

CCES

Assessing d.d.

Paradox

Lessons



Menu

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

" Trio" Identity

Trio

LLP

What's Big

CCES

Assessing d.d.

Paradox

Lessons

Mean Squared Error (MSE)

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

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト ・ ヨ

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



Menu

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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): $D_I = E_R(\hat{\rho}^2)$

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

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト ・ ヨ



Menu

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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): $D_I = E_R(\hat{\rho}^2)$

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

- 日本 - 1 日本 - 日本 - 日本

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



Menu

Xiao-Li Meng Department of Statistics, Harvard University

7

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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 $N^{-1/2}$...



A Law of Large Populations (LLP)





A Law of Large Populations (LLP)

Wenn	
The function	

8

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Bi

Assessing d

Paradox

Lessons



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

Actual Error Benchmark SRS Standard Error

The (lack-of) design effect (Deff)

$${
m Deff} = {{
m MSE}\over {
m Benchmark~SRS~MSE}} = (N-1)D_I$$



A Law of Large Populations (LLP)

Wenn	
The function	

8

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big? CCES Assessing d.d Paradox If $\rho = \mathsf{E}_{\mathcal{R}}(\hat{\rho}) \neq 0$, then on average, the relative error $\uparrow \sqrt{N}$: $\frac{\text{Actual Error}}{\frac{1}{N} - 1} = \sqrt{N - 1}\hat{\rho}$

Benchmark SRS Standard Error

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

$$Deff = \frac{MSE}{Benchmark SRS MSE} = (N-1)D_{I}$$





Effective Sample Size

Menu	9
Xiao-Li M	leng
Statistic	:s, d
What's Bi	g?
CCES	

< □ > < □ > < Ξ > < Ξ > < Ξ > ○ < ♡ < ♡



Effective Sample Size

Menu

Xiao-Li Meng Department of Statistics, Harvard University

9

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big?

CCES

Assessing d.d.

Paradox

Lessons

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

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

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

▲ロト ▲帰 ト ▲ ヨ ト ▲ ヨ ト ● ● ● の Q ()



Effective Sample Size

Menu

Xiao-Li Meng Department of Statistics, Harvard University

9

Motivation

Soup

" Trio" Identity

Trio

LLP

What's Big?

CCES

Assessing d.d. Paradox

Lessons

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

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

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

What matters is the relative size f = n/N

$$n_{\rm eff} = rac{n}{1+(1-f)[(N-1)D_I-1]} pprox rac{f}{1-f}rac{1}{\hat{
ho}^2}.$$

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶



Gaining 2020 Vision: Assessing the behavioral $\hat{\rho}$ using validated voter counts ($\approx 35,000$)

Menu

10

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

CCES

Assessing d.d.i

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)



Gaining 2020 Vision: Assessing the behavioral $\hat{\rho}$ using validated voter counts ($\approx 35,000$)

Menu

Xiao-Li Meng Department of Statistics, Harvard University

Motivatio

Soup

" Trio"

Trio

LLP

What's Bi

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)

・ロト ・ 厚ト ・ ヨト ・ ヨト



Reasonable predictions for Clinton's Vote Share



Gaining 2020 Vision: Assessing the behavioral $\hat{\rho}$ using validated voter counts ($\approx 35,000$)

Menu

Xiao-Li Meng Department of Statistics, Harvard University

Motivatio

Soup

"Trio" Io

Trio

LLP

What's Big

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)



Reasonable predictions for Clinton's Vote Share Serious underestimation of Trump's Vote Share



Assessing $\hat{\rho}$ using state-level data

Menu

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

CCES

Assessing d.d.i

Paradox

Lessons

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

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



Assessing $\hat{\rho}$ using state-level data

Menu

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

" Trio" Identit

Tric

ПP

What's E

CCES

Assessing d.d.i

Lessons

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

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

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト ・ ヨ





Assessing $\hat{\rho}$ using state-level data

Menu

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

" Trio" Identit

Tric

LLP

What's B

CCES

Assessing d.d.i

Lessons

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

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



◆ロト ◆昼 ≻ ◆ 臣 ト ◆ 臣 ト ○ 臣 - • • ○ � @ ●



Menu	12
------	----

Xiao-Li Meng Department of Statistics, Harvard University

Motivatio

Soup

"Trio" Identity

Trio

LLP

What's Big

CCES

Assessing d.d.i

Paradox

Lessons

• Many (major) survey results published before Nov 8, 2016;



Menu

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

CCES

Assessing d.d.i

Paradox

Lessons

- Many (major) survey results published before Nov 8, 2016;
- Roughly amounts to 1% of eligible voters: $n \approx 2,300,000$;

◆□▶ ◆□▶ ◆□▶ ◆□▶ □ の ○ ○

• Equivalent to 2,300 surveys of 1,000 respondents each.



Menu :

Xiao-Li Meng Department of Statistics, Harvard University

Motivatior

Soup

" Trio" Identity

Trio

LLP

What's Big

CCES

Assessing d.d.i

Lessons

- Many (major) survey results published before Nov 8, 2016;
- Roughly amounts to 1% of eligible voters: $n \approx 2,300,000$;
- Equivalent to 2,300 surveys of 1,000 respondents each.

When $\hat{ ho} = -0.005 = -1/200, D_I = 1/40000$, and hence

$$n_{\rm eff} = rac{f}{1-f} rac{1}{D_I} = rac{1}{99} imes 40000 pprox 404!$$



Menu :

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identit

Trio

LLP

What's Big

CCES

Assessing d.d.i

Paradox

Lessons

- Many (major) survey results published before Nov 8, 2016;
- Roughly amounts to 1% of eligible voters: $n \approx 2,300,000$;
- Equivalent to 2,300 surveys of 1,000 respondents each.

When $\hat{
ho} = -0.005 = -1/200, D_I = 1/40000$, 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 :

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identit

Trio

LLP

What's Big

CCES

Assessing d.d.i

Paradox

Lessons

- Many (major) survey results published before Nov 8, 2016;
- Roughly amounts to 1% of eligible voters: $n \approx 2,300,000$;
- Equivalent to 2,300 surveys of 1,000 respondents each.

When $\hat{
ho} = -0.005 = -1/200, D_I = 1/40000$, 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 $\hat{\rho} = -0.005$.
- Butterfly Effect due to Law of Large Populations (LLP)

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



Visualizing LLP: Actual Coverage for Clinton



イロト イポト イヨト イヨト



Visualizing LLP: Actual Coverage for Trump





The Big Data Paradox:

Menu 15

Xiao-Li Meng Department of Statistics, Harvard University

Motivatior

Soup

"Trio" Identi

Trio

LLP

What's Big

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.

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト ・ ヨ



Menu 16

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

CCES

Assessing d.d.i

Paradox

Lessons

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

▲ロト ▲帰 ト ▲ ヨ ト ▲ ヨ ト ● ● ● の Q ()



Menu 16

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

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.



Menu 16

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Tric

LLP

What's Big

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.



Menu 1

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Tric

LLP

What's Big

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: Classical theory is BIG for "big data", as long as we let it go outside the classical box.



In case you are kind enough to invite me again ...

Menu	17
Xiao-Li M Departme Statisti	leng nt of cs,
Harvar Univers	d ity
CCES	
Lessons	



In case you are kind enough to invite me again ...

Menu 1

Xiao-Li Meng Department of Statistics, Harvard University

Motivation

Soup

"Trio" Identity

Trio

LLP

What's Big

CCES

Assessing d.d.i

Paradox

Lessons

The sequel: Meng (2018/9)

Statistical Paradises and Paradoxes in Big Data (II): Multi-resolution Inference, Simpson's Paradox, and Individualized Treatments