

HARVARD DATA SCIENCE REVIEW

A Telescopic, Microscopic, and Kaleidoscopic View of Data Science

HOME JOURNAL 🔻 CATEGORIES 🔻 TOPICS 🗶 MEDIA FEATURES 🗶 PODCAST SUBMISSIONS 👻 ABOUT 👻 MASTHEAD 💌



2021 Best New Journal in Science, Technology, & Medicine

As an open access platform of the Harvard Data Science Initiative, the *Harvard Data Science Review* features **foundational thinking, research milestones, educational innovations, and major applications,** with a primary emphasis on reproducibility, replicability, and readability. It aims to publish contents that help to define and shape data science as a scientifically rigorous and globally impactful multidisciplinary field based on the principled and purposed production, processing, parsing and analysis of data. By uniting the strengths of a premier research journal, a cutting-edge educational publication, and a popular magazine, HDSR provides a crossroads at which fundamental data science research and education intersect directly with societally-important applications from industry, governments, NGOs, and others. By disseminating inspiring, informative, and intriguing articles and media materials, HDSR aspires to be a global forum on *everything data science and data science for everyone*.

ര

 \times



Type the word 'scientist' into your favorite browser and search for images. Most likely you will see photos of actual scientists from various fields. Now repeat the search, using 'data scientist.' You will see far fewer photos, but many animated figures standing by or pointing to various lists of must-have skills that read like tiger parents' assignments for their children.

As we advance deeper into the digital age, our societal demands for data scientists naturally rise in both quantity and quality. Most of us have some knowledge about other kinds of scientists (granted that such knowledge can be quite flawed), but we are much less clear about who data scientists are and what they do. Indeed, what exactly is *data science* (DS)? As you may have guessed, the answer depends on whom you ask. Some say DS is CS (computer science). Others think DS is simply S (statistics). You may even run into someone who declares DS is just hyped-up BS (and I don't mean "Bayesian statistics").

Generation-Collection-Processing-Storage-Management-Analysis-Visualization-Interpretation

Data Conceptualization





Data governance is key to interpretation: The Data Life Cycle Reconceptualizing data in data science

Reconceptualizing data in data science

Sabina Leonelli Sociology, Philosophy and Anthropology, University of Exeter Jeannette M.Wing Computer Science Columbia University

Data Curation and Provenance



The Lives and After Lives of Data

Christine L. Borgman Information Studies, UCLA



Data Science for Governing and Policy Making



An Interview With Murray Edelman on the History of the Exit Poll

by Murray Edelman, Liberty Vittert, and Xiao-Li Meng

An interview with Murray Edelman by Liberty Vittert and Xiao-Li Meng



ABOUT 🔻

Understanding the Quality of the 2020 Census Is Essential

by John Thompson



Recombination: A Family of Markov Chains for Redistricting

by Daryl DeFord, Moon Duchin, and Justin Solomon



The Politics of COVID-19: Partisan Polarization About the Pandemic Has Increased, but Support for Health Care Reform Hasn't Moved at All

by John Sides, Chris Tausanovitch, and Lynn Vavreck



SCRAM: A Platform for Securely Measuring Cyber Risk

by Leo de Castro, Andrew W. Lo, Taylor Reynolds, Fransisca Susan, Vinod Vaikuntanathan, Daniel Weitzner, and Nicolas Zhang



Error, Uncertainty, and the Shifting Ground of Census Data by Dan Bouk



Implementing Differential Privacy: Seven Lessons From the 2020 United States Census by Michael B. Hawes



Measuring the Gross Domestic Product (GDP): The Ultimate Data Science Project by Brian C. Moyer and Abe Dunn





Search Dashboard - Login or Signup

OME ISSUES *

SECTIONS COLUMNS

COLLECTIONS PODCAST SUBMIT ABOU

ABOUT * MASTHEAD *

[™] [™]20[™]/11/10

Differential Privacy for the 2020 U.S. Census: Can We Make Data Both Private and Useful?

Special Issue 2

FROM THE EDITORS



Harnessing the Known Unknowns: Differential Privacy and the 2020 Census

by Ruobin Gong, Erica L. Groshen, and Salil Vadhan

Published: Jun 24, 2022

Special Issue 2: Differential Privacy for the 2020 U.S. Census

CENSUS: IMPORTANCE, HISTORY, AND TECHNICAL CHANGES



Coming to Our Census: How Soci 🗙 🕂		
edu/pub/1g1cbvkv/release/5?readingCollection=a13	33a0a2	
ibe M Gmail 📫 https://hdsr.mitpres 🔀 Inbox - :	xImeng@g 🔀 Your discussion arti 📀 Data Science for th M Revision on the Dat 🌍 outloo	ok 🥠 Oscars 2021: The c 🔓 Google 📧 Top Chinese Song
HD	Search Dash	board 🔻 Login or Signup
НОМЕ	ISSUES • SECTIONS • COLUMNS • COLLECTIONS • PODCAST SUBMIT • ABOUT • MASTHEAD •	0 y 🗴
	Data Science for Governing and Policy Mak ···· 3 more Published on Jan 31, 2020 DOI 10.1162/99608f92.c871f9e0	SHOW DETAILS
	Coming to Our Census: How Social	стте [#] SOCIAL <
	Statistics Underpin Our Democracy (and	DOWNLDAD
	Republic)	
	by Teresa A. Sullivan Published on Jan 31, 2020	
	D. 7 months ago	

ABSTRACT

The 2020 Census provides the opportunity to reflect on the key role statisticians, demographers, and other social scientists play in safeguarding American democracy. Democracy requires numbers for its proper functioning, and there is now a large statistical infrastructure of which the constitutionally mandated census is the keystone. Mistrust of the government is a major obstacle for the census, potentially affecting both accuracy and completeness. The mistrust is stimulated by fears of individual or household census data being willingly or inadvertently shared with other government agencies (data privacy issues) or even foreign actors (hacking). As two 2019 Supreme Court decisions in juxtaposition suggest, no checks or balances protect the integrity of the census. The professional integrity of statisticians is the best defense of the census.

Keywords: census, democracy, statistical infrastructure, data errors, data privacy

This article is accompanied by multiple invited discussion pieces and a rejoinder by the author.





How my own research has been affected/evolved ...

- Data Quality (Multi-source Inference/Learning)
- Data Privacy (Multi-phase Inference/Learning)
- Data Granularity (Multi-resolution Inference/Learning)

Current Issue · Special Issue 3: Personalized (N-of-1) Trials: Methods, Applications, and Impact

See also our special issue, "Differential Privacy fo	r the 2020 U.S. Cens	us: Can We Make Data	Both Private an	d Useful?'	
	Click to Read				

FROM THE EDITORS



Introducing Data Sciences to N-of-1 Designs, Statistics, Use-Cases, the Future, and the Moniker 'N-of-1' Trial

by Karina Davidson, Ken Cheung, Ciaran Friel, and Jerry Suls

Published: Sep 08, 2022

Special Issue 3: Personalized (N-of-1) Trials: Methods, Applications, and Impact

EVERYTHING DATA SCIENCE & DATA SCIENCE FOR EVERYONE

INAUGURAL VOLUME 2019 HDSR.MITPRESS.MIT.EDU

HARVARD DATA SCIENCE REVIEW

What does the educated citizen need to know about Data Science?

What is a data life cycle?

Are we in an Al revolution?

Is there a data creators' advantage?

Who wrote "In My Life"?

What is the Evidence Act?



A Conversation With Steve Ballmer Building Representative Miniatures out of Non-representative Big Data: An Interplay of Data Quantity, Quality, and Privacy

Xiao-Li Meng, Harvard University



∃ → (∃ →

Building Representative Miniatures out of Non-representative Big Data: An Interplay of Data Quantity, Quality, and Privacy

Xiao-Li Meng, Harvard University

- Meng (2022). Miniaturizing Data Defect Correlation: A Versatile Strategy for Handling Non-Probability Samples. *Survey Methodology*.
- Meng (2018) Statistical Paradises and Paradoxes in Big Data (I): From The Law of Large Populations to The Big Data Paradox. *Annals of Applied Statistics*, 12, 685–726.
- Bradley, Sejdinovic, Meng, Kuriwaki, Isakov, Flaxman (2021)
 Unrepresentative Big Surveys Significantly Overestimated COVID
 Vaccination in the US. Nature Dec; 600(7890):695-700.

< ロト < 同ト < ヨト < ヨト

It has been a decade ...

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



It has been a decade ...

- 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)



It has been a decade ...

- 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)



4 ∃ > < ∃ >

• Law of Large Numbers:

Jakob Bernoulli (1713)



イロト イボト イヨト イヨト

• Law of Large Numbers: Jakob Bernoulli (1713)

• Central Limit Theorem: Abraham de Moivre (1733): error $\propto \frac{1}{\sqrt{n}}$: *n* - sample size



.

- Law of Large Numbers: Jakob Bernoulli (1713)
- Central Limit Theorem: Abraham de Moivre (1733): error $\propto \frac{1}{\sqrt{n}}$: *n* - sample size
- Survey Sampling:
 - Graunt (1662); Laplace (1882)



4 B K 4 B K

- Law of Large Numbers: Jakob Bernoulli (1713)
- Central Limit Theorem: Abraham de Moivre (1733): error $\propto \frac{1}{\sqrt{n}}$: *n* - sample size
- Survey Sampling:
 - Graunt (1662); Laplace (1882)
 - The "intellectually violent revolution" in 1895 by Anders Kiær, Statistics Norway

- Law of Large Numbers: Jakob Bernoulli (1713)
- Central Limit Theorem: Abraham de Moivre (1733): error $\propto \frac{1}{\sqrt{n}}$: *n* - sample size
- Survey Sampling:
 - Graunt (1662); Laplace (1882)
 - The "intellectually violent revolution" in 1895 by Anders Kiær, Statistics Norway





- Law of Large Numbers: Jakob Bernoulli (1713)
- Central Limit Theorem: Abraham de Moivre (1733): error $\propto \frac{1}{\sqrt{n}}$: *n* - sample size
- Survey Sampling:
 - Graunt (1662); Laplace (1882)
 - The "intellectually violent revolution" in 1895 by Anders Kiær, Statistics Norway





- Law of Large Numbers: Jakob Bernoulli (1713)
- Central Limit Theorem: Abraham de Moivre (1733): error $\propto \frac{1}{\sqrt{n}}$: *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)



- Law of Large Numbers: Jakob Bernoulli (1713)
- Central Limit Theorem: Abraham de Moivre (1733): error $\propto \frac{1}{\sqrt{n}}$: *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)





Xiao-Li Meng, Harvard University

- Law of Large Numbers: Jakob Bernoulli (1713)
- Central Limit Theorem: Abraham de Moivre (1733): error $\propto \frac{1}{\sqrt{n}}$: *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







Triumphs: Ensuring learning is feasible and cost effective



Triumphs: Ensuring learning is feasible and cost effective

- Census is not feasible in most studies
- Sampling makes a study feasible and even cost effective
- Sampling also reduces privacy cost since it is a random suppression



Triumphs: Ensuring learning is feasible and cost effective

- Census is not feasible in most studies
- Sampling makes a study feasible and even cost effective
- Sampling also reduces privacy cost since it is a random suppression

Troubles: Ensuring "repetitiveness" is increasingly challenging



Triumphs: Ensuring learning is feasible and cost effective

- Census is not feasible in most studies
- Sampling makes a study feasible and even cost effective
- Sampling also reduces privacy cost since it is a random suppression

Troubles: Ensuring "repetitiveness" is increasingly challenging

- Non-response/coverage bias (not merely the rate) is on the rise
- Social media data are non-representative by design
- Administrative data are not probability samples

Thoughts: Handling data defects and privacy concomitantly

Triumphs: Ensuring learning is feasible and cost effective

- Census is not feasible in most studies
- Sampling makes a study feasible and even cost effective
- Sampling also reduces privacy cost since it is a random suppression

Troubles: Ensuring "repetitiveness" is increasingly challenging

- Non-response/coverage bias (not merely the rate) is on the rise
- Social media data are non-representative by design
- Administrative data are not probability samples

Thoughts: Handling data defects and privacy concomitantly

- Can we utilize the inherent data defects as data privacy protections?
- Can we engineer "representative miniatures" that trade data utility and privacy more sensibly than adding noises to defective data?



(日) (四) (三) (三)

Xiao-Li Meng, Harvard University

Recording (Reporting/Responding) Indicator *R* for a **finite population** $\mathbb{X} = \{X_1, \dots, X_N\}$:

 $R_i = 1$, if X_i is included in a sample of size *n*, and $R_i = 0$ otherwise



イロト イボト イヨト イヨト

Recording (Reporting/Responding) Indicator R for a **finite population** $\mathbb{X} = \{X_1, \dots, X_N\}$:

 $R_i = 1$, if X_i is included in a sample of size n, and $R_i = 0$ otherwise

Design Prob: Prob controlled and implemented by human

Ex: Prob sampling, $Pr(R_i = 1 | \mathbb{X}) = n/N$ More: Bootstraps, clinical trials, differential privacy, permutation test, MC.



イロト イボト イヨト イヨト

Recording (Reporting/Responding) Indicator R for a **finite population** $\mathbb{X} = \{X_1, \dots, X_N\}$:

 $R_i = 1$, if X_i is included in a sample of size n, and $R_i = 0$ otherwise

Design Prob: Prob controlled and implemented by human

Ex: Prob sampling, $Pr(R_i = 1 | \mathbb{X}) = n/N$ More: Bootstraps, clinical trials, differential privacy, permutation test, MC.

Divine Prob: Existential Prob by God/Nature or by faith.

Ex: Real-World Evidence, but we assume R is random, and $Pr(R_i = 1|X) = \pi(X_i)$. [Needed for define "missing at Random"]



< □ > < □ > < □ > < □ > < □ >

Recording (Reporting/Responding) Indicator R for a **finite population** $\mathbb{X} = \{X_1, \dots, X_N\}$:

 $R_i = 1$, if X_i is included in a sample of size n, and $R_i = 0$ otherwise

Design Prob: Prob controlled and implemented by human

Ex: Prob sampling, $Pr(R_i = 1 | \mathbb{X}) = n/N$ More: Bootstraps, clinical trials, differential privacy, permutation test, MC.

Divine Prob: Existential Prob by God/Nature or by faith.

Ex: Real-World Evidence, but we assume R is random, and $Pr(R_i = 1|X) = \pi(X_i)$. [Needed for define "missing at Random"]

Device Prob: Probability constructs invoked for analysis.

Ex:
$$\operatorname{logit}[\pi(x)] = \alpha + \beta x$$

More: Prob distributions for expressing belief, prior knowledge, assumptions, idealizations, compromises, desperation.

Xiao-Li Meng, Harvard University

Menu

A Fundamental Identity for Estimation Error (Meng, 2018)



ヨト・イヨト

Image: Image:

Xiao-Li Meng, Harvard University

A Fundamental Identity for Estimation Error (Meng, 2018)

• Population $\{X_1, ..., X_N\}$; Estimand $\bar{\mu}_N = \frac{\sum_{i=1}^N G(X_i)}{N}$;



イロト イボト イヨト イヨト

A Fundamental Identity for Estimation Error (Meng, 2018)

- Population $\{X_1, ..., X_N\}$; Estimand $\bar{\mu}_N = \frac{\sum_{i=1}^N G(X_i)}{N}$;
- Estimator: sample average

$$\bar{\mu}_{n} = \frac{\sum_{i=1}^{N} R_{i} G(X_{i})}{\sum_{i=1}^{N} R_{i}} \equiv \frac{\sum_{i=1}^{N} R_{i} G(X_{i})}{n_{R}}$$

where $R_i = 1$ if X_i is recorded, and zero otherwise.

• Neither X nor R is assumed random.


A Fundamental Identity for Estimation Error (Meng, 2018)

- Population $\{X_1, ..., X_N\}$; Estimand $\bar{\mu}_N = \frac{\sum_{i=1}^N G(X_i)}{N}$;
- Estimator: sample average

$$\bar{\mu}_{n} = \frac{\sum_{i=1}^{N} R_{i} G(X_{i})}{\sum_{i=1}^{N} R_{i}} \equiv \frac{\sum_{i=1}^{N} R_{i} G(X_{i})}{n_{R}}$$

where $R_i = 1$ if X_i is recorded, and zero otherwise.

• Neither X nor R is assumed random.

Expressing the exact error w.r.t $I \sim U[1, ..., N]$ (a divine probability):

$$\bar{\mu}_n - \bar{\mu}_N = \frac{\mathsf{E}_I[R_I G_I]}{\mathsf{E}_I[R_I]} - \mathsf{E}_I[G_I] = \frac{\mathsf{Cov}_I(R_I, G_I)}{\mathsf{E}_I[R_I]}$$
$$= \hat{\rho}_{R,G} \times \sqrt{\frac{N-n}{n}} \times \sigma_G$$

イロト 不得下 イヨト イヨト 二日

A Fundamental Identity for Estimation Error (Meng, 2018)

- Population $\{X_1, ..., X_N\}$; Estimand $\bar{\mu}_N = \frac{\sum_{i=1}^N G(X_i)}{N}$;
- Estimator: sample average

$$\bar{\mu}_{n} = \frac{\sum_{i=1}^{N} R_{i} G(X_{i})}{\sum_{i=1}^{N} R_{i}} \equiv \frac{\sum_{i=1}^{N} R_{i} G(X_{i})}{n_{R}}$$

where $R_i = 1$ if X_i is recorded, and zero otherwise.

• Neither X nor R is assumed random.

Expressing the exact error w.r.t $I \sim U[1, ..., N]$ (a divine probability):

$$\bar{\mu}_n - \bar{\mu}_N = \frac{\mathsf{E}_I[R_I G_I]}{\mathsf{E}_I[R_I]} - \mathsf{E}_I[G_I] = \frac{\mathsf{Cov}_I(R_I, G_I)}{\mathsf{E}_I[R_I]}$$
$$= \hat{\rho}_{R,G} \times \sqrt{\frac{N-n}{n}} \times \sigma_G$$

イロト イポト イヨト イヨト

Why teaching is so important ...



イロト イポト イヨト イヨト

Xiao-Li Meng, Harvard University

Why teaching is so important ...





7 / 26

Xiao-Li Meng, Harvard University

Menu

2 December 1993

Statistics & Probability Letters 18 (1993) 345-348 North-Holland

On the absolute bias ratio of ratio estimators

Xiao-Li Meng Department of Statistics, University of Chicago, IL, USA

Received January 1993 Revised March 1993

Abstract: The elegant Hartley-Ross inequality on the absolute bias ratio (ABR = [Bia] (XE, D) of a ordinary ratio estimator is birst generalized to that is a sportar taral estimator with stratified sampling. It is shown that, as long as the numerators and denominators used to form strata ratios are unbiased estimators, the absolute bias ratio of a separate ratio estimator will neve exceed the square root of the sum of squares of the coefficient of variation of the denominators scrass strata. This provides, at design stages, a simple bound in practice to assess the limit and magnitude of the bias ratio of any separate ratio estimator that shares the same demoninators. East expressions for biases of separate ratio estimators are also given.

Keywords: Combined ratio estimator; separate ratio estimator; stratified sampling.

1. Biases of ordinary ratio estimators

In sample surveys, ordinary ratio estimators are typically employed to estimate (i) a population total, Y, (ii) a population mean, \overline{Y} , or (iii) a population ratio, Y/X. In all of these cases, the ratio estimator has the form

$$r = \frac{\bar{y}}{\bar{x}}Q, \qquad (1.1)$$

where \bar{x} and \bar{y} are the sample means of variable x and y, respectively, and Q is a known quantity. In cases (i) and (ii), Q is the population total and mean of variable x, X and \bar{X} respectively, and the ratio estimator r of (1.1) is used to increase the precision in estimating Y and \bar{Y} by taking

advantage of the positive correlation between yand x and the known values of X and \overline{X} in the population. In case (iii), Q = 1, and the population quantities of variable x need not be known. A comprehensive treatment of ratio estimator (1.1) and other variations can be found in Cochran (1977, Chapter 6).

It is well known that in general, r of (11) is biased for $R = \langle \overline{Y}, \overline{X} \rangle Q$. However, this bias is typically unimportant because it is negligible compared to the standard error of r. An elementary but elegant proof of this fact was given in Hartley and Ross (1954), who noticed the following simple identity

$$E(r) - \frac{E(\bar{y})}{E(\bar{x})}Q = -\frac{\operatorname{Cov}(r, \bar{x})}{E(\bar{x})}.$$
 (1)

.2)

< ∃⇒



 $\underline{\bar{\mu}_n - \bar{\mu}_N} = \hat{\rho}_{R,G} \quad \times$ Exact Error Data Quality



イロト イボト イヨト イヨト





★ Ξ ► < Ξ ►</p>





A B + A B +

< A[™]



• Nothing is assumed random, only using divine probability via 1.



4 B K 4 B K



- Nothing is assumed random, only using divine probability via 1.
- Only assumption: sampled X_i are in the target population.





- Nothing is assumed random, only using divine probability via 1.
- Only assumption: sampled X_i are in the target population.

Compare to the additive "Variance+Bias" decomposition

$$\underbrace{\bar{\mu}_n - \bar{\mu}_N}_{\text{Exact Error}} = \underbrace{\bar{\mu}_n - \mathsf{E}_R(\bar{\mu}_n)}_{\text{Sampling Error}} + \underbrace{\mathsf{E}_R(\bar{\mu}_n) - \bar{\mu}_N}_{\text{Bias}}$$





- Nothing is assumed random, only using divine probability via 1.
- Only assumption: sampled X_i are in the target population.

Compare to the additive "Variance+Bias" decomposition

$$\underbrace{\bar{\mu}_n - \bar{\mu}_N}_{\text{Exact Error}} = \underbrace{\bar{\mu}_n - \mathsf{E}_R(\bar{\mu}_n)}_{\text{Sampling Error}} + \underbrace{\mathsf{E}_R(\bar{\mu}_n) - \bar{\mu}_N}_{\text{Bias}}$$

- Must involve prob on *R*, design, divine, or device.
- It's "phenotype" decomposition, not "genotype."



• Think about tasting soup



(B)

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





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





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







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









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



(a) < (a)





• But what happens when we fail to stir (well)



Probability Sampling miniaturizes the ddc

For SRS, because $V_{R}(Z) = 1$, where

$$Z = \frac{\bar{\mu}_n - \bar{\mu}_N}{\sqrt{\frac{1-f}{n}}S_G} = \hat{\rho}_{R,G}\sqrt{N-1} \quad \Rightarrow \quad \mathcal{V}_R(\hat{\rho}_{R,G}) = \frac{1}{N-1}$$

Hence $\hat{\rho}_{R,G} \in \left(-\frac{3}{\sqrt{N-1}}, \frac{3}{\sqrt{N-1}}\right)$ 99% of the time regardless of G .



(B)

Probability Sampling miniaturizes the *ddc*

For SRS, because $V_{R}(Z) = 1$, where

$$Z = \frac{\bar{\mu}_n - \bar{\mu}_N}{\sqrt{\frac{1-f}{n}}S_G} = \hat{\rho}_{R,G}\sqrt{N-1} \quad \Rightarrow \quad \mathcal{V}_R(\hat{\rho}_{R,G}) = \frac{1}{N-1}$$

Hence $\hat{\rho}_{R,G} \in \left(-\frac{3}{\sqrt{N-1}}, \frac{3}{\sqrt{N-1}}\right)$ 99% of the time regardless of *G*.

$$\bar{\mu}_{n} - \bar{\mu}_{N} = \underbrace{\hat{\rho}_{R,G}}_{1/\sqrt{N}} \times \underbrace{\sqrt{\frac{N-n}{n}}}_{\sqrt{N}} \times \sigma_{G}$$



Probability Sampling miniaturizes the *ddc*

For SRS, because $V_{_R}(Z) = 1$, where

$$Z = \frac{\bar{\mu}_n - \bar{\mu}_N}{\sqrt{\frac{1-f}{n}}S_G} = \hat{\rho}_{R,G}\sqrt{N-1} \quad \Rightarrow \quad \mathcal{V}_R(\hat{\rho}_{R,G}) = \frac{1}{N-1}$$

Hence $\hat{\rho}_{R,G} \in \left(-\frac{3}{\sqrt{N-1}}, \frac{3}{\sqrt{N-1}}\right)$ 99% of the time regardless of *G*.

$$\bar{\mu}_{n} - \bar{\mu}_{N} = \underbrace{\hat{\rho}_{R,G}}_{1/\sqrt{N}} \times \underbrace{\sqrt{\frac{N-n}{n}}}_{\sqrt{N}} \times \sigma_{G}$$

• This cancellation is **THE reason** that we can ignore *N* for any *G*

ddc:
$$\hat{\rho}_{R,G}$$
$$\frac{\bar{\mu}_n - \bar{\mu}_N}{\sqrt{\frac{1-f}{n}}S_G} = \hat{\rho}_{R,G}\sqrt{N-1}$$

ddi:
$$D_I = \mathsf{E}_R(\hat{\rho}_{R,G}^2)$$

Deff $\equiv \frac{\mathrm{MSE}(\bar{\mu}_n)}{\mathrm{SRS MSE}} = D_I(N-1)$



ddc:
$$\hat{\rho}_{R,G}$$

$$\frac{\bar{\mu}_n - \bar{\mu}_N}{\sqrt{\frac{1-f}{n}}S_G} = \hat{\rho}_{R,G}\sqrt{N-1}$$

$$Deff \equiv \frac{\text{MSE}(\bar{\mu}_n)}{\text{SRS MSE}} = D_I(N-1)$$

 $ddi = \frac{Deff}{N-1} =$ Design effect per subject in the population

• For SRS:
$$D_I = (N-1)^{-1}$$



ddc:
$$\hat{\rho}_{R,G}$$

 $\frac{\bar{\mu}_n - \bar{\mu}_N}{\sqrt{\frac{1-f}{n}}S_G} = \hat{\rho}_{R,G}\sqrt{N-1}$

$$ddi: D_I = \mathsf{E}_R(\hat{\rho}_{R,G}^2)$$

$$\mathrm{Deff} \equiv \frac{\mathrm{MSE}(\bar{\mu}_n)}{\mathrm{SRS MSE}} = D_I(N-1)$$

 $ddi = \frac{Deff}{N-1} =$ Design effect per subject in the population

• For SRS:
$$D_I = (N - 1)^{-1}$$

• Probability sample $\implies D_I \propto N^{-1}$



ddc:
$$\hat{\rho}_{R,G}$$

$$\frac{\bar{\mu}_n - \bar{\mu}_N}{\sqrt{\frac{1-f}{n}}S_G} = \hat{\rho}_{R,G}\sqrt{N-1}$$
Deff $\equiv \frac{\text{MSE}(\bar{\mu}_n)}{\text{SRS MSE}} = D_I(N-1)$

 $ddi = \frac{Deff}{N-1} =$ Design effect per subject in the population

• For SRS:
$$D_I = (N-1)^{-1}$$

• Probability sample \implies $D_I \propto N^{-1}$

•
$$MSE(\bar{\mu}_n) \propto n^{-1} \iff D_I \propto N^{-1}$$

ddc:
$$\hat{\rho}_{R,G}$$

$$\frac{\bar{\mu}_n - \bar{\mu}_N}{\sqrt{\frac{1-f}{n}}S_G} = \hat{\rho}_{R,G}\sqrt{N-1}$$
Deff $\equiv \frac{\text{MSE}(\bar{\mu}_n)}{\text{SRS MSE}} = D_I(N-1)$

 $ddi = \frac{Deff}{N-1} =$ **Design effect per subject in the population**

• For SRS:
$$D_I = (N-1)^{-1}$$

- Probability sample $\implies D_I \propto N^{-1}$
- $MSE(\bar{\mu}_n) \propto n^{-1} \iff D_I \propto N^{-1}$

Deep Trouble

- when D_I does not vanish with N^{-1} ;
- or equivalently when $\hat{\rho}_{\rm \scriptscriptstyle G,R}$ does not vanish with $N^{-1/2}$...

э

・ロト ・ 日 ト ・ ヨ ト ・



э

(日) (四) (三) (三)

Xiao-Li Meng, Harvard University

The Effective Sample Size of a "Big Data" in terms of SRS size $n_{\text{eff}} = \frac{n}{1 + (1 - f)[(N - 1)\mathsf{E}_R[\hat{\rho}_{R,G}^2] - 1]} \approx \frac{f}{1 - f} \frac{1}{\hat{\rho}^2}$



< ロト < 同ト < ヨト < ヨト

The Effective Sample Size of a "Big Data" in terms of SRS size $n_{\text{eff}} = \frac{n}{1 + (1 - f)[(N - 1)\mathsf{E}_R[\hat{\rho}_{R,G}^2] - 1]} \approx \frac{f}{1 - f} \frac{1}{\hat{\rho}^2}$

Why do we need random testing? (Walter Dempsey, Twitter)

• NY State: $N \approx 20$ M;



The Effective Sample Size of a "Big Data" in terms of SRS size $n_{\text{eff}} = \frac{n}{1 + (1 - f)[(N - 1)\mathsf{E}_R[\hat{\rho}_{R,G}^2] - 1]} \approx \frac{f}{1 - f} \frac{1}{\hat{\rho}^2}$

Why do we need random testing? (Walter Dempsey, Twitter)

- NY State: $N \approx 20$ M;
- Suppose we conduct n = 10,000 COVID test: f = 1/2000



The Effective Sample Size of a "Big Data" in terms of SRS size $n_{\text{eff}} = \frac{n}{1 + (1 - f)[(N - 1)\mathsf{E}_R[\hat{\rho}^2_{R,G}] - 1]} \approx \frac{f}{1 - f} \frac{1}{\hat{\rho}^2}$

Why do we need random testing? (Walter Dempsey, Twitter)

- NY State: $N \approx 20$ M;
- Suppose we conduct n = 10,000 COVID test: f = 1/2000
- Suppose the selective testing resulted in $\hat{\rho} = 0.005$;

$$n_{
m eff} = rac{0.0005}{0.9995} imes rac{1}{0.005^2} pprox 20$$



The Effective Sample Size of a "Big Data" in terms of SRS size $n_{\text{eff}} = \frac{n}{1 + (1 - f)[(N - 1)\mathsf{E}_R[\hat{\rho}^2_{R,G}] - 1]} \approx \frac{f}{1 - f} \frac{1}{\hat{\rho}^2}$

Why do we need random testing? (Walter Dempsey, Twitter)

- NY State: $N \approx 20$ M;
- Suppose we conduct n = 10,000 COVID test: f = 1/2000
- Suppose the selective testing resulted in $\hat{\rho} = 0.005$;

$$n_{\rm eff} = \frac{0.0005}{0.9995} \times \frac{1}{0.005^2} \approx 20$$

• Hence $\hat{\rho} = 0.005$ implies a 99.80% loss of sample size!

Where Did $\hat{\rho} = 0.005$ Come From? (Meng, 2018)

Cooperative Congressional Election Study by Ansolabehere, Schaffner, Luks, Rivers on Oct 4 - Nov 6, 2016 (YouGov); Analysis assisted by Shiro Kuriwaki



イロト イボト イヨト イヨト

Where Did $\hat{\rho} = 0.005$ Come From? (Meng, 2018)

Cooperative Congressional Election Study by Ansolabehere, Schaffner, Luks, Rivers on Oct 4 - Nov 6, 2016 (YouGov); Analysis assisted by Shiro Kuriwaki





< ロト < 同ト < ヨト < ヨト



Where Did $\hat{\rho} = 0.005$ Come From? (Meng, 2018)

Cooperative Congressional Election Study by Ansolabehere, Schaffner, Luks, Rivers on Oct 4 - Nov 6, 2016 (YouGov); Analysis assisted by Shiro Kuriwaki



Studies of COVID-19 vaccine uptake vary in design

	Axios-Ipsos Coronavirus Index	Census Household Pulse	Facebook-Delphi COVID Symptom Survey
Mode	online	online	online
Sampling	lpsos	Census Master	Facebook Active User
frame	KnowledgePanel	Address File	Base
n	1,000/wave	65,000/wave	250,000/wave
Question Wording	"Do you personally know anyone who has already received the COVID-19 vaccine?" Answers include "Yes, I have received the vaccine"	"Have you received a COVID-19 vaccine?"	"Have you had a COVID-19 vaccination?"
Target	2019 CPS March	2018 ACS, 1-year	2018 CPS March
Population	Supplement, US 18+	est., US 18+	Supplement, US 18+
Weighting Variables	gender × age, race, education, Census region, metropolitan status, household income, partisanship	metropolitan statistical area (MSA), state x education x gender x age, state x Hispanic ethnicity x gender x age	state x age x gender and "proprietary covariates"
Big Data Paradox: The Bigger the Data, The Surer We Fool Ourselves



Estimated Effective Sample Sizes (dropping three 0s!)



Dramatic Reduction in Effective Sample Size (> 99.9%)



Composition of Survey Respondents

	Composition of U.S. Adults							Survey Estimates		
	Ax	ios-Ipsos	Household Pulse		Delphi-Facebook		ACS	Household Pulse		
Education	Raw	Weighted	Raw	Weighted	Raw	Weighted	Benchmark	Vax	Will	Hes
High School	35%	39%	14%	39%	19%	21%	39%	39%	40%	21%
Some College	29	30	32	30	36	36	30	44	38	18
4-Year College	19	17	29	17	25	25	19	54	36	10
Post-Graduate	17	14	26	13	20	18	11	67	26	7

	Composition of U.S. Adults							Survey Estimates		
	Axios-Ipsos		Household Pulse		Delphi-Facebook		ACS	Household Pulse		
Race/Ethnicity	Raw	Weighted	Raw	Weighted	Raw	Weighted	Benchmark	Vax	Will	Hes
White	71%	63%	75%	62%	74%	68%	60%	50%	33%	17%
Black	10	12	7	11	6	6	12	42	39	19
Hispanic	11	16	10	17	11	16	16	38	48	14
Asian			5	5	2	3	6	51	43	5



イロト イポト イヨト イヨト

"Weight, Weight, Don't tell me ..."

There are always those "messy" weights ...

$$\bar{\mu}_{w} = \frac{\sum_{i=1}^{N} R_{i} w_{i} G(X_{i})}{\sum_{i=1}^{N} R_{i} w_{i}}$$



A B + A B +

"Weight, Weight, Don't tell me ..."

There are always those "messy" weights ...

$$\bar{\mu}_w = \frac{\sum_{i=1}^N R_i w_i G(X_i)}{\sum_{i=1}^N R_i w_i}$$

Let CV_w be the coefficient of variation of W_l given $R_l = 1$

$$\bar{\mu}_{w} - \bar{\mu}_{N} = \hat{\rho}_{\mathrm{Rw},\mathrm{G}} \times \sqrt{\frac{1 - f_{w}}{f_{w}}} \times \sigma_{\mathrm{G}}$$

where

$$f_w = \frac{n_w}{N}, \qquad n_w = \frac{n}{1 + CV_w^2}$$
 (Kish, 1965)



4 ∃ > < ∃ >

"Weight, Weight, Don't tell me ..."

There are always those "messy" weights ...

$$\bar{\mu}_w = \frac{\sum_{i=1}^N R_i w_i G(X_i)}{\sum_{i=1}^N R_i w_i}$$

Let CV_{W} be the coefficient of variation of W_{I} given $R_{I} = 1$

$$\bar{\mu}_{\rm W} - \bar{\mu}_{\rm N} = \hat{\rho}_{\rm Rw,G} \times \sqrt{\frac{1 - f_{\rm W}}{f_{\rm W}}} \times \sigma_{\rm G}$$

where

$$f_w = \frac{n_w}{N}, \qquad n_w = \frac{n}{1 + CV_w^2}$$
 (Kish, 1965)

• Seeking w to make $\hat{\rho}_{_{Rw},_G} < \hat{\rho}_{_{R,G}}$ and compensate for $n_w < n$.



Estimation Methods for Non-probability Samples

- A non-probability sample $\{(y_i, X_i), i \in S\}, S = \{i : R_i = 1\}$
- An auxiliary probability sample $\{X_i, i \in S^*\}, \qquad S^* = \{i : R^*_i = 1\}$
- Key Assumption: $y_i \perp R_i | X_i$ (Missing at Random)
- A device model: p(y, R|x) = p(y|x)p(R|x)



イロト イボト イヨト イヨト

Estimation Methods for Non-probability Samples

- A non-probability sample $\{(y_i, X_i), i \in S\}, S = \{i : R_i = 1\}$
- An auxiliary probability sample $\{X_i, i \in S^*\}, \qquad S^* = \{i: R^*_i = 1\}$
- Key Assumption: $y_i \perp R_i | X_i$ (Missing at Random)
- A device model: p(y, R|x) = p(y|x)p(R|x)

Quasi-randomization: Estimate $\pi(x) = \Pr_p(R = 1|x)$

$$\hat{\mu} = \frac{\sum_{i=1}^{N} R_i w_i y_i}{\sum_{i=1}^{N} R_i w_i}, \qquad w_i \propto \hat{\pi}^{-1}(X_i)$$



< ロト < 同ト < ヨト < ヨト

Estimation Methods for Non-probability Samples

- A non-probability sample $\{(y_i, X_i), i \in S\}, S = \{i : R_i = 1\}$
- An auxiliary probability sample $\{X_i, i \in S^*\}, \qquad S^* = \{i: R^*_i = 1\}$
- Key Assumption: $y_i \perp R_i | X_i$ (Missing at Random)
- A device model: p(y, R|x) = p(y|x)p(R|x)

Quasi-randomization: Estimate $\pi(x) = \Pr_{p}(R = 1|x)$

$$\hat{\mu} = \frac{\sum_{i=1}^{N} R_i w_i y_i}{\sum_{i=1}^{N} R_i w_i}, \qquad w_i \propto \hat{\pi}^{-1}(X_i)$$

Super-population/model assisted: also fit y = m(x)

$$\hat{\mu}_{+} = \frac{\sum_{i=1}^{N} R_{i} w_{i} (y_{i} - \hat{m}(X_{i}))}{\sum_{i=1}^{N} R_{i} w_{i}} + \frac{\sum_{i=1}^{N} R_{i}^{*} \hat{m}(X_{i})}{\sum_{i=1}^{N} R_{i}^{*}},$$

< 150 ►

()

$$\hat{\mu}_{+} - \bar{y}_{N} = \frac{\mathsf{Cov}_{I}(R_{I}w_{I}, y_{I} - \hat{m}(X_{I}))}{\mathsf{E}_{I}(R_{I}w_{I})} + \frac{\mathsf{Cov}_{I}(R_{I}^{*}, \hat{m}(X_{I}))}{\mathsf{E}_{I}(R_{I}^{*})},$$



Э

・ロト ・ 四ト ・ ヨト ・ ヨト

22 / 26

$$\hat{\mu}_{+} - \bar{y}_{N} = \frac{\operatorname{Cov}_{i}(R_{i}w_{i}, y_{i} - \hat{m}(X_{i}))}{\mathsf{E}_{i}(R_{i}w_{i})} + \frac{\operatorname{Cov}_{i}(R_{i}^{*}, \hat{m}(X_{i}))}{\mathsf{E}_{i}(R_{i}^{*})},$$

Under a divine model for (R, Y|X) and sampling model for R^*

 $\mathsf{E}[\hat{\mu}_{+}] - \mu \sim \mathsf{E}_{\mathsf{x}} \left\{ \mathsf{Cov}_{\mathsf{I}}[\pi_{\mathsf{I}} \mathsf{w}_{\mathsf{I}}, \delta_{\mathsf{I}}] \right\}, \quad \text{where } \delta_{\mathsf{I}} \equiv \mathsf{E}(\mathsf{y}_{\mathsf{I}}|\mathsf{X}_{\mathsf{I}}) - \hat{m}(\mathsf{X}_{\mathsf{I}})$



イロト イボト イヨト イヨト

$$\hat{\mu}_{+} - \bar{y}_{N} = \frac{\operatorname{Cov}_{I}(R_{I}w_{I}, y_{I} - \hat{m}(X_{I}))}{\mathsf{E}_{I}(R_{I}w_{I})} + \frac{\operatorname{Cov}_{I}(R_{I}^{*}, \hat{m}(X_{I}))}{\mathsf{E}_{I}(R_{I}^{*})},$$

Under a divine model for (R, Y|X) and sampling model for R^*

 $\mathsf{E}[\hat{\mu}_{+}] - \mu \sim \mathsf{E}_{\mathsf{x}} \left\{ \mathsf{Cov}_{\mathsf{I}}[\pi_{\mathsf{I}} \mathsf{w}_{\mathsf{I}}, \delta_{\mathsf{I}}] \right\}, \quad \text{where } \delta_{\mathsf{I}} \equiv \mathsf{E}(\mathsf{y}_{\mathsf{I}}|\mathsf{X}_{\mathsf{I}}) - \hat{m}(\mathsf{X}_{\mathsf{I}})$

• Quasi-randomization: making $\pi_I w_I \propto 1$, (Q)



イロト イポト イヨト イヨト

$$\hat{\mu}_{+} - \bar{y}_{N} = \frac{\operatorname{Cov}_{i}(R_{i}w_{i}, y_{i} - \hat{m}(X_{i}))}{\mathsf{E}_{i}(R_{i}w_{i})} + \frac{\operatorname{Cov}_{i}(R_{i}^{*}, \hat{m}(X_{i}))}{\mathsf{E}_{i}(R_{i}^{*})},$$

Under a divine model for (R, Y|X) and sampling model for R^*

 $\mathsf{E}[\hat{\mu}_{+}] - \mu \sim \mathsf{E}_{\mathsf{x}} \left\{ \mathsf{Cov}_{\mathsf{I}}[\pi_{\mathsf{I}} \mathsf{w}_{\mathsf{I}}, \delta_{\mathsf{I}}] \right\}, \quad \text{where } \delta_{\mathsf{I}} \equiv \mathsf{E}(\mathsf{y}_{\mathsf{I}}|\mathsf{X}_{\mathsf{I}}) - \hat{m}(\mathsf{X}_{\mathsf{I}})$

- Quasi-randomization: making $\pi_I w_I \propto 1$, (Q)
- Super-population: making $\delta_I \equiv E(y_I|X_I) \hat{m}(X_I) = 0,$ (S)



イロト イポト イヨト イヨト

$$\hat{\mu}_{+} - \bar{y}_{N} = \frac{\operatorname{Cov}_{I}(R_{I}w_{I}, y_{I} - \hat{m}(X_{I}))}{\mathsf{E}_{I}(R_{I}w_{I})} + \frac{\operatorname{Cov}_{I}(R_{I}^{*}, \hat{m}(X_{I}))}{\mathsf{E}_{I}(R_{I}^{*})},$$

Under a divine model for (R, Y|X) and sampling model for R^*

 $\mathsf{E}[\hat{\mu}_{+}] - \mu \sim \mathsf{E}_{\mathsf{x}} \left\{ \mathsf{Cov}_{\mathsf{I}}[\pi_{\mathsf{I}} \mathsf{w}_{\mathsf{I}}, \delta_{\mathsf{I}}] \right\}, \quad \text{where } \delta_{\mathsf{I}} \equiv \mathsf{E}(\mathsf{y}_{\mathsf{I}}|\mathsf{X}_{\mathsf{I}}) - \hat{m}(\mathsf{X}_{\mathsf{I}})$

- Quasi-randomization: making $\pi_I w_I \propto 1$, (Q)
- Super-population: making $\delta_I \equiv E(y_I|X_I) \hat{m}(X_I) = 0,$ (S)
- Doubly robust: either (Q) or (S) makes Cov(π_iw_i, δ_i) = 0, but we don't need to know which one.



イロト イヨト イヨト

$$\hat{\mu}_{+} - \bar{y}_{N} = \frac{\mathsf{Cov}_{I}(R_{I}w_{I}, y_{I} - \hat{m}(X_{I}))}{\mathsf{E}_{I}(R_{I}w_{I})} + \frac{\mathsf{Cov}_{I}(R_{I}^{*}, \hat{m}(X_{I}))}{\mathsf{E}_{I}(R_{I}^{*})},$$

Under a divine model for (R, Y|X) and sampling model for R^*

 $\mathsf{E}[\hat{\mu}_{+}] - \mu \sim \mathsf{E}_{\mathsf{x}} \left\{ \mathsf{Cov}_{\mathsf{I}}[\pi_{\mathsf{I}} \mathsf{w}_{\mathsf{I}}, \delta_{\mathsf{I}}] \right\}, \quad \text{where } \delta_{\mathsf{I}} \equiv \mathsf{E}(\mathsf{y}_{\mathsf{I}}|\mathsf{X}_{\mathsf{I}}) - \hat{m}(\mathsf{X}_{\mathsf{I}})$

- Quasi-randomization: making $\pi_I w_I \propto 1$, (Q)
- Super-population: making $\delta_I \equiv E(y_I|X_I) \hat{m}(X_I) = 0,$ (S)
- Doubly robust: either (Q) or (S) makes Cov(π_iw_i, δ_i) = 0, but we don't need to know which one.
- "Double+ robustness": the validity holds if and only if

$$\mathsf{E}_{\mathsf{x}}\left\{\mathsf{Cov}_{\mathsf{y}}\left[\pi_{\mathsf{y}}\mathsf{w}_{\mathsf{y}},\delta_{\mathsf{y}}\right]\right\}=\mathsf{0}$$



イロト イヨト イヨト イヨト

• If $\Pr(R_i = 1 | X, Y) = \pi_i$, sub-sample with $\Pr(S_i = 1 | R_i = 1) \propto \pi_i^{-1}$



イロト イポト イヨト イヨト

- If $\Pr(R_i = 1 | X, Y) = \pi_i$, sub-sample with $\Pr(S_i = 1 | R_i = 1) \propto \pi_i^{-1}$
- Hence $Pr(S_iR_i = 1) \propto 1$, creating an equal-probability sample



イロト イボト イヨト イヨト

- If $\Pr(R_i = 1 | X, Y) = \pi_i$, sub-sample with $\Pr(S_i = 1 | R_i = 1) \propto \pi_i^{-1}$
- Hence $Pr(S_iR_i = 1) \propto 1$, creating an equal-probability sample

The Effective Sample Size of "Big Data" in terms of SRS size

$$n_B pprox rac{f_B}{1-f_B}rac{1}{\hat{
ho}_B^2},$$

where $f_B = n/N$ is the *relative size* of sample B.



- If $\Pr(R_i = 1 | X, Y) = \pi_i$, sub-sample with $\Pr(S_i = 1 | R_i = 1) \propto \pi_i^{-1}$
- Hence $\Pr(S_i R_i = 1) \propto 1$, creating an equal-probability sample

The Effective Sample Size of "Big Data" in terms of SRS size

$$n_B pprox rac{f_B}{1-f_B}rac{1}{\hat{
ho}_B^2},$$

where $f_B = n/N$ is the *relative size* of sample B.

Trading quantity for quality: create an unweighted sub-sample D such that

$$\frac{f_B}{1-f_B}\frac{1}{\hat{\rho}_B^2} < \frac{f_D}{1-f_D}\frac{1}{\hat{\rho}_D^2},$$

where $f_D = f_B f_S$, and f_S is the sub-sampling rate.



< ロト < 同ト < ヨト < ヨト

23 / 26

Estimand p = P(Y = 1); estimable $p^* = P(Y = 1|R = 1)$



イロト イボト イヨト イヨト

Estimand p = P(Y = 1); estimable $p^* = P(Y = 1|R = 1)$

• Reporting rates:

$$r_y = \Pr(R = 1 | Y = y); \quad r = \frac{r_1}{r_0}$$

• Sub-sampling rates:

$$s_y = \Pr(S = 1 | R = 1, Y = y);$$
 $s = \frac{s_1}{s_0}$



イロト イポト イヨト イヨト

Estimand p = P(Y = 1); estimable $p^* = P(Y = 1|R = 1)$

Reporting rates:

$$r_y = \Pr(R = 1 | Y = y); \quad r = \frac{r_1}{r_0}$$

Sub-sampling rates:

,

$$s_y = \Pr(S = 1 | R = 1, Y = y);$$
 $s = \frac{s_1}{s_0}$

Counterbalancing: (r-1)(s-1) < 0

(i) If
$$\delta = r - 1 > 0$$
, then take any s

$$rac{[1-(1-p^*)\delta]_+}{1+(1+p^*)\delta} \le s < 1$$
 (1)

(ii) If $\delta = r - 1 < 0$, then take any s

$$1 < s \le \frac{1 - (1 - p^*)\delta}{[1 + (1 + p^*)\delta]_+}$$
(2)





イロト イボト イヨト イヨト



We do not know r, but

• Suppose a previous survey had r = 1.5

Statistics

< ロト < 同ト < ヨト < ヨト



We do not know r, but

- Suppose a previous survey had r = 1.5
- We might feel comfortable to assume that the current $r \in (1.2, 1.8)$

Statistic



We do not know r, but

- Suppose a previous survey had r = 1.5
- We might feel comfortable to assume that the current $r \in (1.2, 1.8)$
- Suppose we observe *p*^{*} = 0.6

Statistic



We do not know r, but

- Suppose a previous survey had r = 1.5
- We might feel comfortable to assume that the current $r \in (1.2, 1.8)$
- Suppose we observe *p*^{*} = 0.6
- Then the max of the lower bound is 0.7

Statistic



We do not know r, but

- Suppose a previous survey had r = 1.5
- We might feel comfortable to assume that the current $r \in (1.2, 1.8)$
- Suppose we observe $p^* = 0.6$
- Then the max of the lower bound is 0.7
- If r = 1.5, then s = 1/1.5 = 0.67 will be optimal, but any $s \in [0.7, 1)$ will lead to smaller MSE compared to not using CBS

Statistics

イロト イヨト イヨト イヨト

There is no free lunch

• Data quality and data quantity trade-off



There is no free lunch

- Data quality and data quantity trade-off
- Data cleanness and data relevance trade-off



There is no free lunch

- Data quality and data quantity trade-off
- Data cleanness and data relevance trade-off
- Data utility and data privacy trade-off



There is no free lunch

. . .

- Data quality and data quantity trade-off
- Data cleanness and data relevance trade-off
- Data utility and data privacy trade-off



There is no free lunch

. . .

- Data quality and data quantity trade-off
- Data cleanness and data relevance trade-off
- Data utility and data privacy trade-off



There is no free lunch

- Data quality and data quantity trade-off
- Data cleanness and data relevance trade-off
- Data utility and data privacy trade-off

• ...

But let's not overpay either

- Avoid add more noise than necessary
- Avoid weighted data when an unweighted sub-sample can provide similar statistical information

