

# **Survey Research and Design**

## **Survey Sampling and Mode**

---

William Marble

September 19, 2023

## Discussion

Rank preferences for discussion – link on Canvas

Problem set 1 distributed today after class → start brainstorming survey topic  
+ meet with me to discuss

R tutoring: Monday 12 - 1:30 p.m. and Tuesday 1 - 3 p.m.

## **Sampling with Unequal Inclusion Probabilities**

---

## Survey Sampling

- ▶ Last class: simple random sampling

## Survey Sampling

- ▶ Last class: simple random sampling
- ▶ Each member of the sampling frame has an equal chance of selection

## Survey Sampling

- ▶ Last class: simple random sampling
- ▶ Each member of the sampling frame has an equal chance of selection
- ▶ Easy analysis: unweighted mean is unbiased

## Unequal Inclusion Probabilities

Some designs have unequal inclusion probabilities:

- ▶ Cluster or stratified sampling with unequal cluster sizes
- ▶ Oversamples of populations of interest

## Unequal Inclusion Probabilities

Some designs have unequal inclusion probabilities:

- ▶ Cluster or stratified sampling with unequal cluster sizes
- ▶ Oversamples of populations of interest

Unweighted means will no longer be unbiased  $\leadsto$  need to account for design decisions in analysis

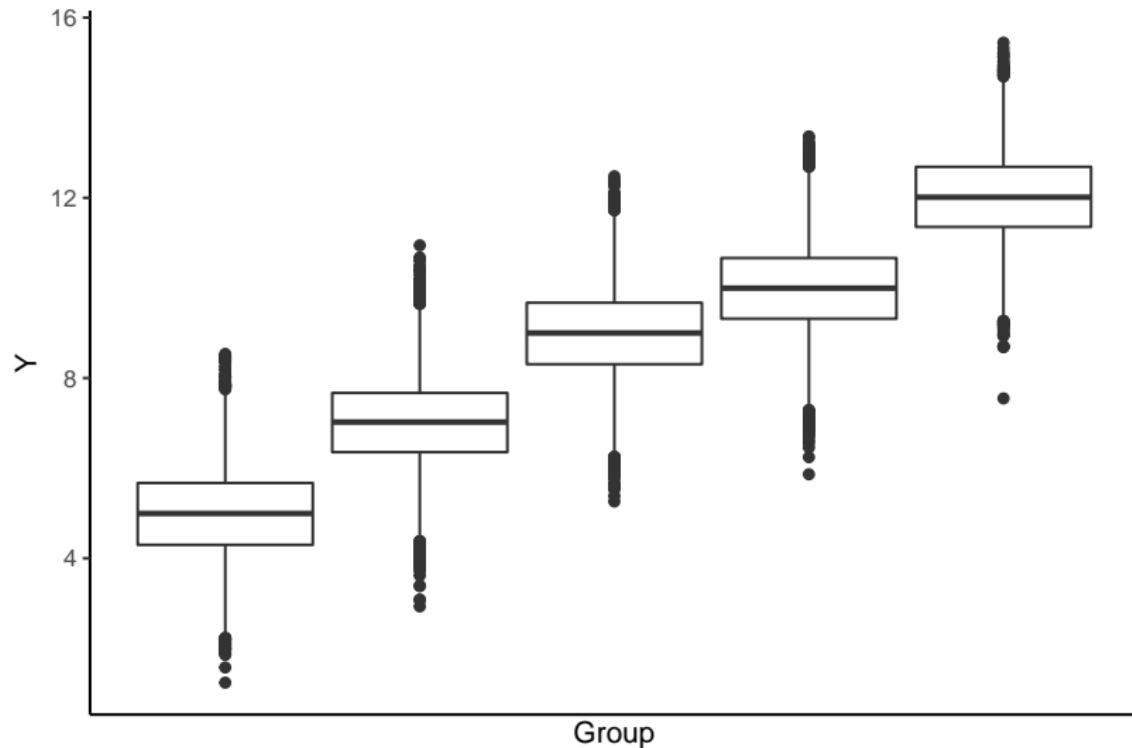
## Unequal Inclusion Probabilities

Consider stratified sampling where we select  $n$  units from each strata, which each have population sizes  $N_k$ . What is the probability that any given unit is included?

(See illustration in R)

## Data Simulation

There are 5 strata in the population, and people in different strata have different means:



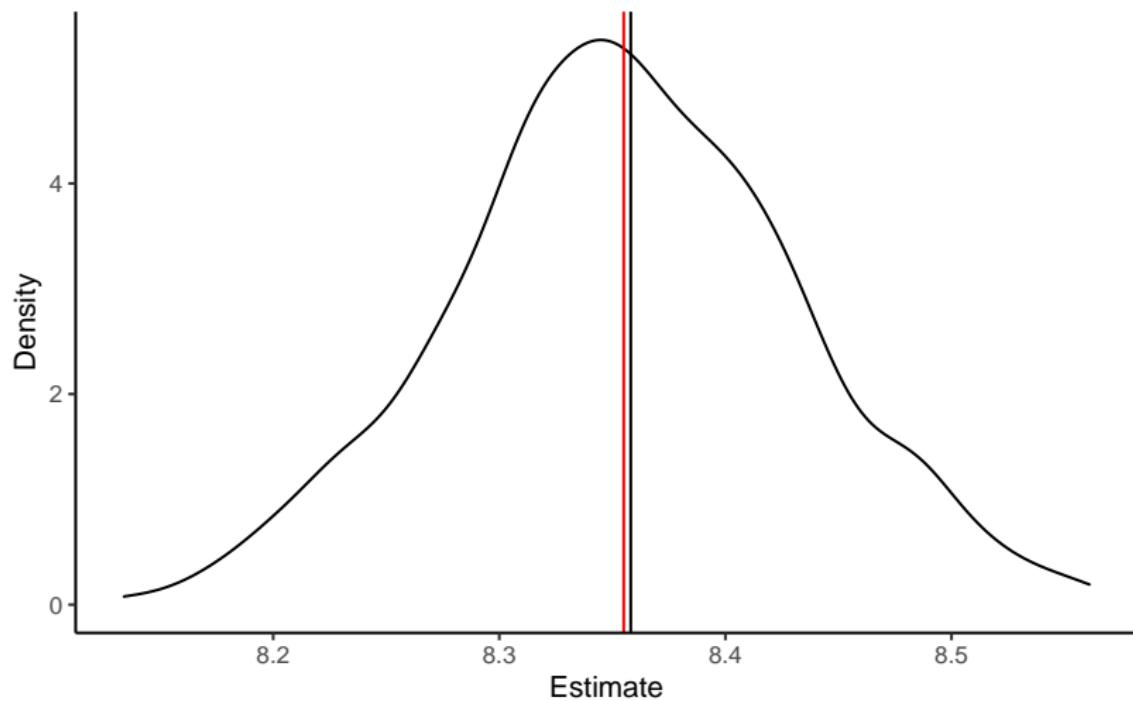
## Data Simulation

Cluster	Proportion of Population	Mean
1	0.19	5
2	0.24	7
3	0.21	9
4	0.24	10
5	0.12	12
<b>Total</b>	1	8.36

## Simple Random Sampling

Under SRS, we get an unbiased estimate of the mean, without weighting

Simple Random Sampling is Unbiased

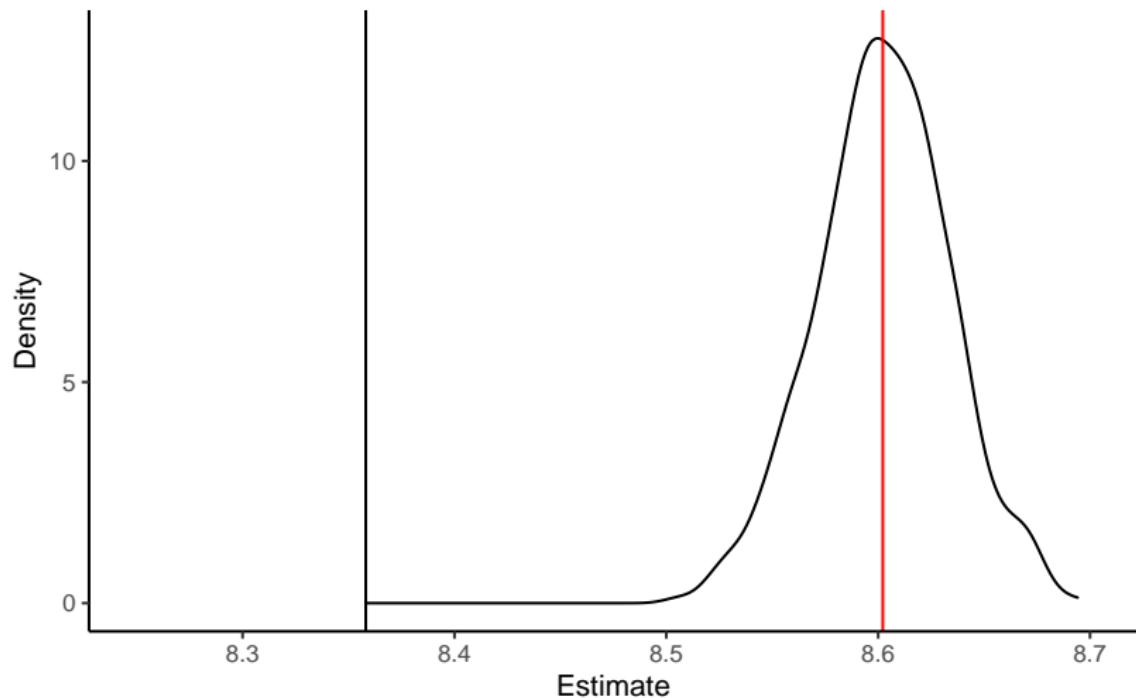


## Stratified Sampling

Instead, imagine we sample 200 people from each cluster (i.e. stratified random sampling). Will the sample mean be unbiased for the population mean? Why or why not?

## Sampling Distribution of Stratified Samples

Stratified Sampling Without Weight Adjustment is Biased



## Weighting by Inverse Probability of Sampling

- ▶ Probability any unit gets sampled is  $p_i = 200/N_{k[i]}$ , where  $N_{k[i]}$  is size of strata for unit  $i$
- ▶ Samples will over-represent small clusters — smaller  $N_{k[i]}$  means higher  $p$

## Weighting by Inverse Probability of Sampling

- ▶ Probability any unit gets sampled is  $p_i = 200/N_{k[i]}$ , where  $N_{k[i]}$  is size of strata for unit  $i$
- ▶ Samples will over-represent small clusters — smaller  $N_{k[i]}$  means higher  $p$
- ▶ Solution: take weighted mean where weight  $w_i = 1/p_i$
- ▶ Higher weight for under-represented strata

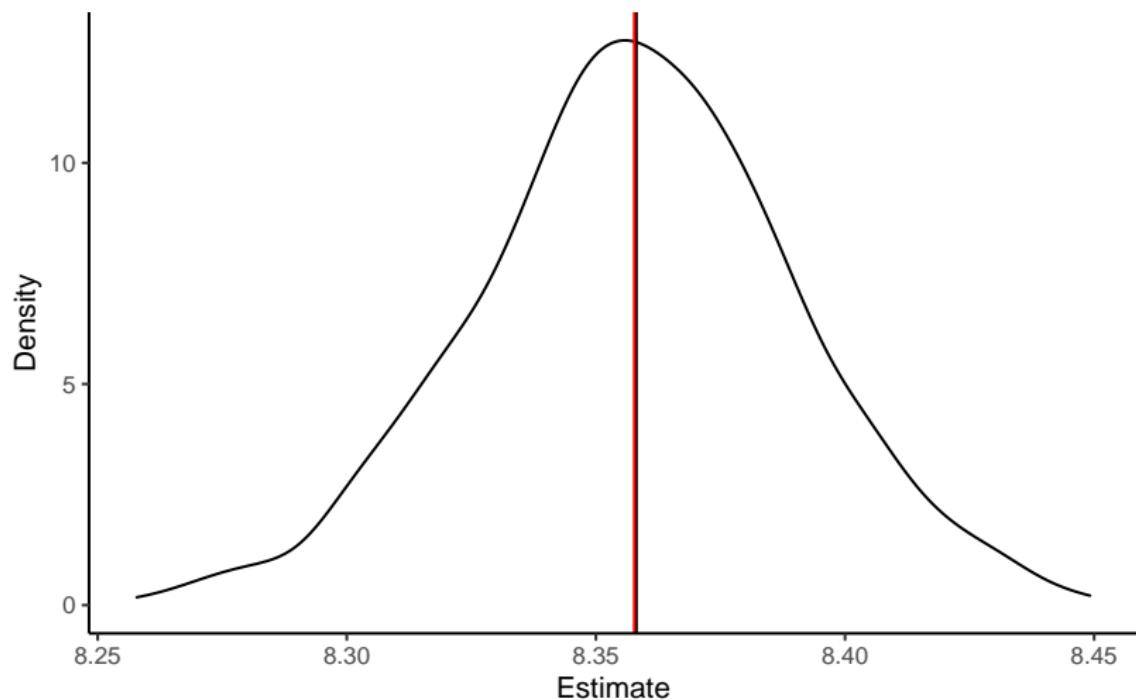
## Weighting by Inverse Probability of Sampling

- ▶ Probability any unit gets sampled is  $p_i = 200/N_{k[i]}$ , where  $N_{k[i]}$  is size of strata for unit  $i$
- ▶ Samples will over-represent small clusters — smaller  $N_{k[i]}$  means higher  $p$
- ▶ Solution: take weighted mean where weight  $w_i = 1/p_i$
- ▶ Higher weight for under-represented strata

In R: `weighted.mean(y, w = 1/p)`

## Sampling Distribution of Stratified Samples

Stratified Sampling With Weight Adjustment is Unbiased



## Survey Weighting, Generally

- ▶ Weighting accounts for differences between the population and sample
- ▶ Weight by inverse probability of sampling when sample inclusion probs are unequal (can get complicated!)
- ▶ Need adjustments to variance estimation too

## Survey Weighting, Generally

- ▶ Weighting accounts for differences between the population and sample
- ▶ Weight by inverse probability of sampling when sample inclusion probs are unequal (can get complicated!)
- ▶ Need adjustments to variance estimation too
- ▶ For more complex issues (e.g. nonresponse), need additional assumptions beyond knowledge of sampling design
- ▶ Intuition is the same: upweight the under-represented, downweight the over-represented

## **Survey Mode**

---

## Survey Mode

- ▶ So far we've glossed over how to recruit respondents
- ▶ Many ways to recruit respondents, with advantages and disadvantages
- ▶ There are advantages and disadvantages to each method

## **Random Digit Dialing**

## **Random Digit Dialing**

### Advantages

## **Random Digit Dialing**

### Advantages

- ▶ Every phone number has an equal chance of inclusion
- ▶ Cheap to generate sampling frame and sample
- ▶ Can target area codes
- ▶ Almost everyone has a phone

## **Random Digit Dialing**

### Advantages

- ▶ Every phone number has an equal chance of inclusion
- ▶ Cheap to generate sampling frame and sample
- ▶ Can target area codes
- ▶ Almost everyone has a phone

### Downsides

## Random Digit Dialing

### Advantages

- ▶ Every phone number has an equal chance of inclusion
- ▶ Cheap to generate sampling frame and sample
- ▶ Can target area codes
- ▶ Almost everyone has a phone

### Downsides

- ▶ Spam calls are out of control  $\leadsto$  people don't answer
- ▶ People move: incorrect area code
- ▶ Double-coverage for people with multiple phone numbers
- ▶ Hard to target subpopulations

# Postcard Recruitment



Center for Opinion Research  
PO Box 3003  
Lancaster, PA 17604

To complete a survey go to:

[fandm.edu/pagovnt](http://fandm.edu/pagovnt)

or, call:

1-866-366-7655

Before

September 8th, 2019

and provide your survey access  
code to begin the survey:

332464



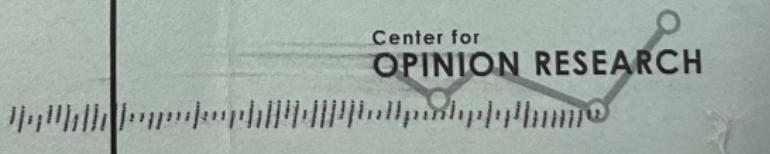
U.S. POSTAGE » PITNEY BOWES



ZIP 17603 \$ 000.35<sup>0</sup>

02 4W 0000336993 AUG 19 2019

Stephen Pettigrew



# Postcard Recruitment

You have been selected from a random sample of registered voters in Pennsylvania to participate in a survey about Pennsylvania state government and issues in politics. The results of the survey will be used to help state legislators gain a better understanding of voters' opinions.

**You may receive a call from 717-358-4632 in Lancaster, PA  
between now and Sunday, September 8th requesting your participation.**

You may also complete the survey online, by visiting the website listed below and entering your survey access code to begin the survey.

**WEBSITE: [fandm.edu/pagovnt](http://fandm.edu/pagovnt)**

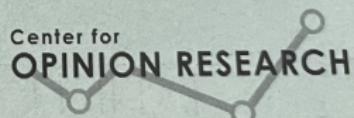
**ACCESS CODE: 332464**

Or, you may call the Center for Opinion Research at **1-866-366-7655** (Monday -Friday 9am - 9pm, Saturday 11am – 3pm, Sunday 3pm – 7pm) and provide your access code listed above.

The survey will take about 10 minutes to complete. Be assured your survey responses will be kept strictly confidential and will only be used for research purposes.

**Thank you for your time and consideration.**

**FRANKLIN & MARSHALL  
COLLEGE**



# **Postcard Recruitment**

## Advantages

## Postcard Recruitment

### Advantages

- ▶ More convenient for respondents
- ▶ Might reach a broader sample than phone or web
- ▶ Can link respondents to geographic location

## Postcard Recruitment

### Advantages

- ▶ More convenient for respondents
- ▶ Might reach a broader sample than phone or web
- ▶ Can link respondents to geographic location

### Disadvantages

## **Postcard Recruitment**

### **Advantages**

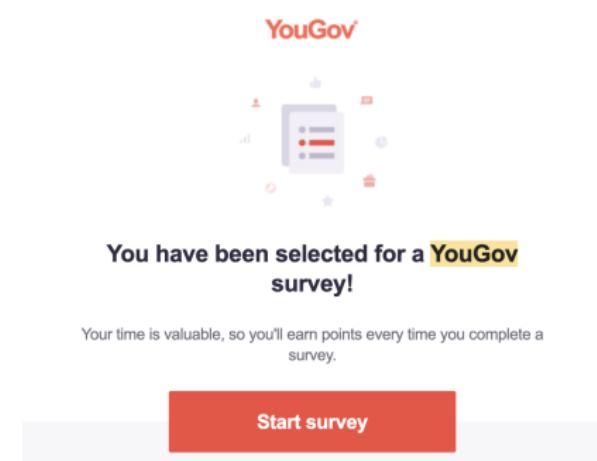
- ▶ More convenient for respondents
- ▶ Might reach a broader sample than phone or web
- ▶ Can link respondents to geographic location

### **Disadvantages**

- ▶ Cost
- ▶ Outdated address info

# Internet Surveys

## Advantages



The image shows a YouGov survey interface. At the top right is the YouGov logo. Below it is a central icon representing a survey, surrounded by various small icons like a person, a lightbulb, and a gear. In the center, the text "You have been selected for a YouGov survey!" is displayed in bold black font. Below this, a smaller text states: "Your time is valuable, so you'll earn points every time you complete a survey." At the bottom right is a red button with the white text "Start survey".

YouGov

You have been selected for a YouGov survey!

Your time is valuable, so you'll earn points every time you complete a survey.

Start survey

# Internet Surveys

## Advantages

- ▶ Many methods of recruitment
- ▶ Fast, convenient, and cheap (!!)
- ▶ Can target niche audiences

YouGov<sup>®</sup>



You have been selected for a YouGov  
survey!

Your time is valuable, so you'll earn points every time you complete a survey.

Start survey

# Internet Surveys

## Advantages

- ▶ Many methods of recruitment
- ▶ Fast, convenient, and cheap (!!)
- ▶ Can target niche audiences

## Disadvantages

YouGov<sup>®</sup>



You have been selected for a YouGov  
survey!

Your time is valuable, so you'll earn points every time you complete a survey.

Start survey

# Internet Surveys

## Advantages

- ▶ Many methods of recruitment
- ▶ Fast, convenient, and cheap (!!)
- ▶ Can target niche audiences

## Disadvantages

- ▶ Usually non-probability samples
- ▶ Unclear who takes online surveys
- ▶ Inattentive responding

YouGov<sup>®</sup>



You have been selected for a YouGov survey!

Your time is valuable, so you'll earn points every time you complete a survey.

**Start survey**

# **Internet Surveys**

## Advantages

## Internet Surveys

### Advantages

- ▶ Many methods of recruitment
- ▶ Fast, convenient, and cheap (!!)
- ▶ Can target niche audiences

## Internet Surveys

### Advantages

- ▶ Many methods of recruitment
- ▶ Fast, convenient, and cheap (!!)
- ▶ Can target niche audiences

### Disadvantages

# Internet Surveys

## Advantages

- ▶ Many methods of recruitment
- ▶ Fast, convenient, and cheap (!!)
- ▶ Can target niche audiences

## Disadvantages

- ▶ Usually non-probability samples
- ▶ Unclear who takes online surveys
- ▶ Inattentive responding



Liverpool F.C. Research

August 30, 2018 · 6

...

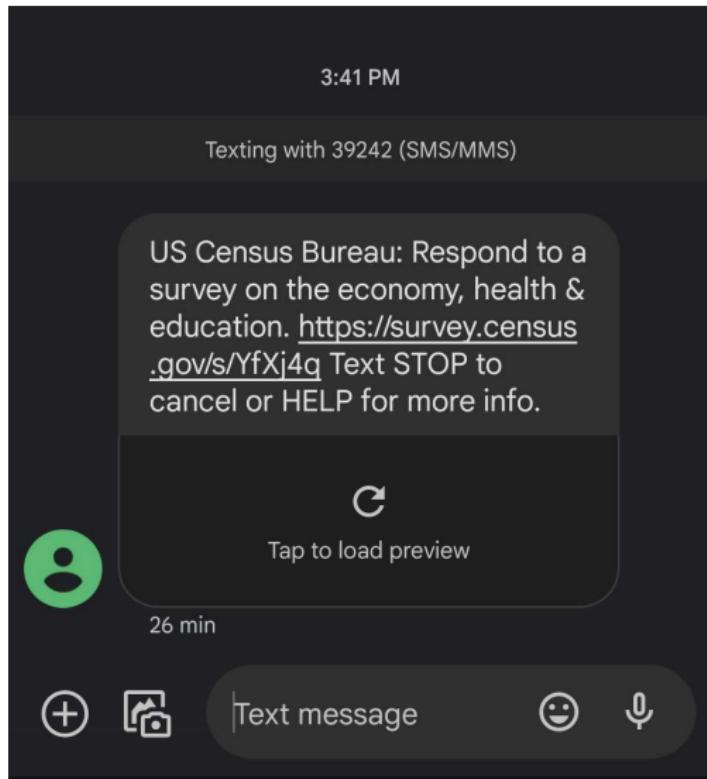
Love the Reds? We are Stanford University researchers studying the views of Liverpool fans on a range of issues. Take this 2 minute survey to help us understand Liverpool supporters! <https://stanforduniversity.qualtrics.com/.../SV...>

STANFORDUNIVERSITY.QUALTRICS.COM

Online Survey Software | Qualtrics Survey Solutions



## Text to Web



## Door-to-Door



©Dave Coverly. All rights reserved.

## Door-to-Door

### Advantages

- ▶ Often used for the highest-quality surveys (e.g. Census, American National Election Studies)
- ▶ Higher-quality responses and more engagement
- ▶ Higher response rates

### Disadvantages

- ▶ Very expensive (!!)
- ▶ Logistically difficult
- ▶ Usually requires more complicated design (and thus analysis)
- ▶ Not easy to follow up with nonrespondents

**Anscombe and Schaffner (2014),  
“Does Survey Mode Still Matter?”**

---

## **Research Question**

What question or problem does this paper address? Why is this question important?

# Research Question

What question or problem does this paper address? Why is this question important?

## 1 Introduction

The rapid increase of Internet penetration in American homes has made web-based polling a viable and affordable alternative for students of public opinion. A particularly affordable method for conducting Internet surveys relies on recruiting volunteers to take Internet polls and then generating representative samples of the target populations of interest from these panels using a technique called sample matching. This is the approach used by YouGov (formerly Polimetrix), which is the firm most commonly used by political scientists conducting research using Internet surveys that utilize opt-in panels.<sup>1</sup> The popularity of this approach has been reflected in the publication of opt-in Internet survey data in the top journals of political science. Between 2006 and 2011, thirty articles published in the *American Political Science Review*, *American Journal of Political Science*, or *Journal of Politics* utilized opt-in Internet survey data produced by YouGov.<sup>2</sup>

Among public opinion researchers, there is a lack of consensus about the acceptability of opt-in Internet surveys. It is difficult to evaluate different methods because we lack evidence that can guide decisions. In 2010, the American Association of Public Opinion Researchers (AAPOR)

## Methods Overview

What methods does the paper propose to use to address the problem?

## Methods Overview

What methods does the paper propose to use to address the problem?

This article presents results from a study conducted in 2010 comparing opt-in Internet, telephone, and mail survey modes. These surveys were designed so that identical questionnaires were used across modes, utilizing questions that are commonly asked on existing surveys (such as those fielded by Pew, the National Health Indicators Survey [NHIS], and ANES), many of which can provide objective baselines against which to compare the estimates.

Our analysis examines three key comparisons. First, we measure the extent to which each approach produces accurate point estimates for measures on which we have validated benchmark data, such as from the US Census. Second, we examine cross-mode similarities and differences in point estimates for political measures (such as attitudes and reported behavior) that cannot be validated. Third, we compare the correlation structure and regression models estimated from data from each mode (e.g., Sanders et al. 2007).

## **Phone Survey**

How did the researchers recruit respondents in the phone survey?

## Phone Survey

How did the researchers recruit respondents in the phone survey?

The telephone survey was conducted January 28–30, with 807 interviews completed with respondents reached via landline numbers and a supplement of one hundred interviews conducted using cell phone numbers. Respondents were selected using random digit dialing, with the cell phone supplement coming from random digit dialing of known cell phone exchanges. Live interviewers were used to administer the questionnaire to respondents. Each telephone number

## **Opt-In Internet Panel**

How did the researchers recruit the online sample?

## Opt-In Internet Panel

How did the researchers recruit the online sample?

- ▶ First recruit panelists:

The Internet sample for our study came from the YouGov online panel. The selection process for this panel includes recruiting a large number of people to serve on the survey panel through various methods, including online advertising. Individuals who join the panel earn rewards (i.e., points that

## Opt-In Internet Panel

How did the researchers recruit the online sample?

- ▶ First recruit panelists:

The Internet sample for our study came from the YouGov online panel. The selection process for this panel includes recruiting a large number of people to serve on the survey panel through various methods, including online advertising. Individuals who join the panel earn rewards (i.e., points that

- ▶ Generate a pseudo-sample from government data and other surveys:

Since YouGov does not use probability sampling to recruit panelists, they instead rely on sample matching to generate representative samples from their panel. When YouGov is commissioned to conduct a survey, they begin by taking a random sample from the target population. For example, if a client is asking for a survey of one thousand American adults, YouGov might draw a random sample from the Census Bureau's American Community Survey (ACS) and use this as the target for constructing a sample from their own panel. In addition to all of the demographic information that is part of the ACS, YouGov is also able to weight on additional factors that it has matched to the ACS. For example, data on reported voter registration and turnout from the Current Population

## Opt-In Internet Panel

How did the researchers recruit the online sample?

- ▶ First recruit panelists:

The Internet sample for our study came from the YouGov online panel. The selection process for this panel includes recruiting a large number of people to serve on the survey panel through various methods, including online advertising. Individuals who join the panel earn rewards (i.e., points that

- ▶ Generate a pseudo-sample from government data and other surveys:

Since YouGov does not use probability sampling to recruit panelists, they instead rely on sample matching to generate representative samples from their panel. When YouGov is commissioned to conduct a survey, they begin by taking a random sample from the target population. For example, if a client is asking for a survey of one thousand American adults, YouGov might draw a random sample from the Census Bureau's American Community Survey (ACS) and use this as the target for constructing a sample from their own panel. In addition to all of the demographic information that is part of the ACS, YouGov is also able to weight on additional factors that it has matched to the ACS. For example, data on reported voter registration and turnout from the Current Population

- ▶ Invite panelists who match pseudo-sample to take the survey:

and ideology were matched from the 2007 Pew US Religious Landscape Survey. Thus, once YouGov draws the target sample from the database, they know what each member of their random sample should look like on a range of demographic, political, and religious characteristics and using these characteristics an algorithm selects the closest matching individuals from their Internet panel to essentially replace each person that was randomly selected into the target sample (Rivers 2006). After matching everyone in the target sample with at least one person

## Mail Survey

How did the researchers recruit the mail sample?

## Mail Survey

How did the researchers recruit the mail sample?

- ▶ Randomly sample addresses from a large list:

The mail survey was generated by mailing questionnaires to 6600 addresses selected randomly from a list provided by a data vendor. The sample was randomly divided into different types of incentive conditions—19% received no incentive, 39% were offered \$1, 39% were offered \$2, and 3% were offered \$5. The overall response rate for this sample was (RR3) 21.1%.

## Mail Survey

How did the researchers recruit the mail sample?

- ▶ Randomly sample addresses from a large list:

The mail survey was generated by mailing questionnaires to 6600 addresses selected randomly from a list provided by a data vendor. The sample was randomly divided into different types of incentive conditions—19% received no incentive, 39% were offered \$1, 39% were offered \$2, and 3% were offered \$5. The overall response rate for this sample was (RR3) 21.1%.

- ▶ Invite them to return the survey by mail or take it online:

Individuals receiving the mail questionnaire were offered the opportunity to either return their survey by post or go online to take the survey. Of those responding to the mail solicitation, 27.5%

## **Analysis**

How do the researchers analyze the data?

## Analysis

How do the researchers analyze the data?

- ▶ Compare responses to benchmarks

The most valuable metric for understanding the validity of a survey mode is to compare it to a validated baseline that can be treated as a valid approximation to the population parameter. Thus, we begin by examining the TSE of each mode by determining the extent to which each survey produces accurate estimates of characteristics for which we have validated benchmarks. Since

## Analysis

How do the researchers analyze the data?

- ▶ Compare responses to benchmarks

The most valuable metric for understanding the validity of a survey mode is to compare it to a validated baseline that can be treated as a valid approximation to the population parameter. Thus, we begin by examining the TSE of each mode by determining the extent to which each survey produces accurate estimates of characteristics for which we have validated benchmarks. Since

- ▶ Several variables: homeownership, cigarette usage, health insurance, turnout and vote choice in 2008

## **Findings on Questions with Baseline Data**

What were the main findings of the paper?

# Findings on Questions with Baseline Data

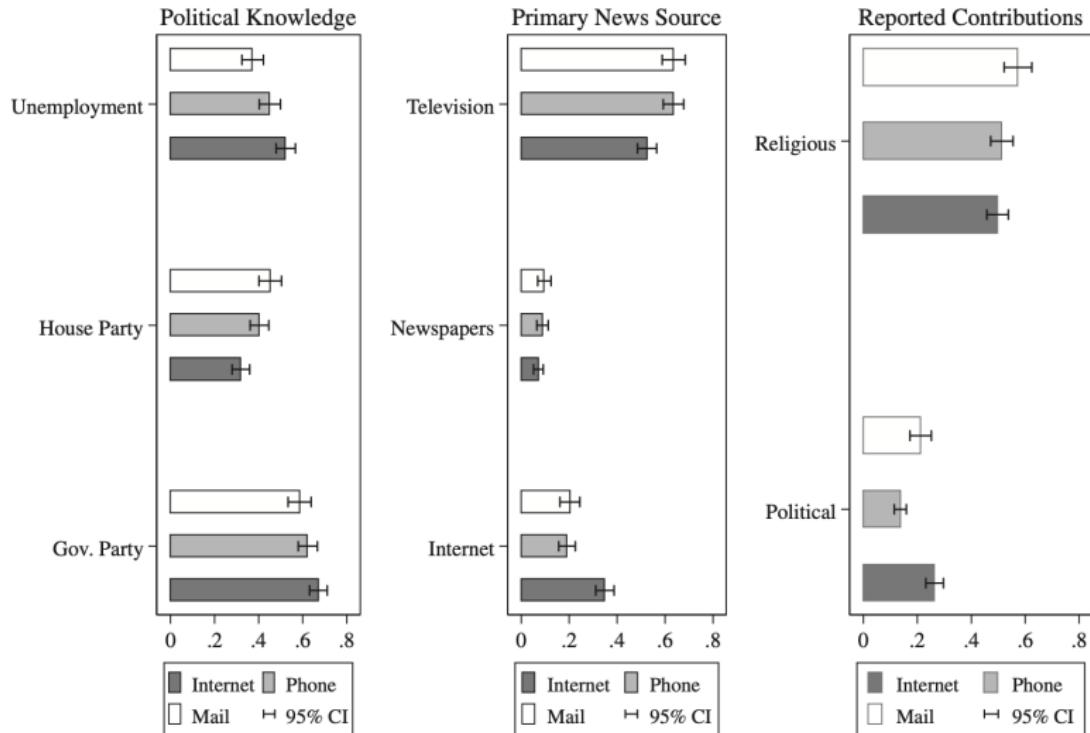
What were the main findings of the paper?

**Table 2** TSE comparison of point estimates by mode

<i>Item</i>	<i>Response</i>	<i>Internet</i>	<i>Phone</i>	<i>Mail</i>	<i>Validating source</i>
Home ownership	Own	<b>0.613</b> (0.573, 0.653)	0.637 (0.584, 0.690)	0.636 (0.594, 0.678)	0.669 (CPS)
Mobility	Moved in past year	0.152 (0.121, 0.183)	0.148 (0.103, 0.192)	0.146 (0.112, 0.180)	0.154 (ACS)
	At address five or more years	0.555 (0.515, 0.595)	0.613 (0.561, 0.665)	0.548 (0.506, 0.590)	0.588 (ACS)
Smoked one hundred cigarettes	Yes	<b>0.504</b> (0.464, 0.544)	<b>0.483</b> (0.432, 0.534)	<b>0.481</b> (0.440, 0.523)	0.430 (NHIS)
Smoke cigarettes now	Every or some days	<b>0.259</b> (0.222, 0.296)	<b>0.248</b> (0.204, 0.293)	0.223 (0.186, 0.261)	0.203 (NHIS)
Health Ins.	None	0.157 (0.128, 0.186)	<b>0.237</b> (0.187, 0.287)	0.154 (0.120, 0.189)	0.167 (SIPP)
Voted in 2008 (if registered)	Yes	0.888 (0.865, 0.911)	0.876 (0.841, 0.911)	<b>0.825</b> (0.788, 0.861)	0.896 (CPS)
Voting method in 2008	By mail	0.191 (0.161, 0.221)	<b>0.122</b> (0.093, 0.150)	0.164 (0.133, 0.194)	0.164 (CPS)
	Early in-person	0.136 (0.111, 0.162)	<b>0.104</b> (0.078, 0.130)	<b>0.105</b> (0.082, 0.128)	0.143 (CPS)
Vote choice in 2008	Obama	<b>0.482</b> (0.444, 0.521)	<b>0.457</b> (0.405, 0.508)	0.553 (0.512, 0.593)	0.529
	McCain	0.474 (0.436, 0.513)	0.502 (0.450, 0.553)	0.432 (0.391, 0.472)	0.456
<b>Average difference</b>		<b>0.031</b>	<b>0.041</b>	<b>0.029</b>	
<b>MSE</b>		<b>0.013</b>	<b>0.021</b>	<b>0.013</b>	

*Note.* SIPP = Survey of Income and Program Participation. Estimates in bold had confidence intervals that did not include the validated figure.

## Findings on Political Knowledge and Interest



**Fig. 2** Comparison of point estimates on knowledge, news source, and reported contributions. Entries are weighted proportions of respondents in each category after excluding those responding “don’t know” or “not sure,” except for the knowledge questions, where “don’t know” or “not sure” was coded as incorrect.

## **Conclusions**

What is the main conclusion of the paper?

## Conclusions

What is the main conclusion of the paper?

Overall, it appears that researchers will not consistently get more accurate results, nor reach substantially different conclusions, when using one mode relative to another. That said, costs are undoubtedly an important consideration for most researchers. The mail mode was particularly expensive, both in terms of actual costs per completed interview and in terms of the extensive time period required to collect an adequate number of responses. The calculation may be a bit closer when

## Discussion

- ▶ Did the researchers leave out any important factors from this study?
- ▶ Do you think the findings would be different if they re-did the study today?
- ▶ How could you improve the study?
- ▶ Any other limitations, questions, concerns, points of confusion or disagreement, etc.?

## Next Week

Writing survey questions

Brainstorm topics for class survey