Summary of Voter Match Research at AAPOR 2018

Edward Paul Johnson, RN-SSI

June 29, 2018

Presentations Summarized

The Truth is Out There: Using Voter Files to Improve Election Polls

Assessing Coverage Bias in Registration-Based Sample (Seth Brohinsky, Abt Associates & Scott Clement, Washington Post)

Examining Coverage and Response Bias Between Data Collection Modes in Voter Files (Edward Paul Johnson, RN-SSI & Nathan LaCombe, The Data Trust)

Examining the Accuracy of Likely Voter Models (Scott Clement, Washington Post)

Relational validity: A new approach to evaluating political surveys (Jonathan Robinson, Catalist & Kevin Collins, ChangePoint Analytics)

Verifying Voter Registration Records (Michael McDonald, University of Florida)

Panel or Wallpaper?: How to Cover your Survey Needs and other advice on...Online Panels

Creating a Probability Framework for Online Panels by Matching to Voter Files (Edward Paul Johnson, RN-SSI & Kori Bishop, Aristotle)



Assessing Coverage Bias in Registration-Based Sample: Overall Research

Problem:

Phone RDD interviewing of voters is getting expensive.

- Decreased Response Rate
- Increased Cell Phone Cost

Voter File Sampling offers solutions:

- Extensive frame information
- Efficiently target registered voter population
- Improved contact rates
- Access to out-of-area cell phones for state and local surveys
- Lower cost per interview

However:

- Missing phone numbers in the voter file can be up to 30%
- Can't use census data to develop reliable weighting targets

What bias is associated with the missing phone numbers?



Research:

Used Aristotle Voter File in Virginia

Found important differences in frame with and without phone numbers

Dual Frame RDD/Voter File Survey

- Two waves (combined N=2,359)
- Tried to match RDD sample back to voter file on first/last name, zip, age.
- Very little matched back to voter file records with no phone numbers (N=96).

Assessing Coverage Bias in Registration-Based Sample: Skews in Aristotle VA Voter File Only with Phone

61% Marta I Status Urk nown 20% 7% 26% 21% Ethnicity Unknown 6% 26% No Phone Number Phone Number Sample Without Phone Less Likely to Vote 42% 17% 70% Vote din 2014 or 2016 General Election 86% 22% 15% No Phone Number Phone Number Sample Without Phone More Republican 10% 28% 12% Phone Number Phone Number No Phone Number No Phone Number 53% ■ 18 - 29 ■ 30 - 39 ■ 40 - 49 ■ 50 - 69 ■ 70 + Demo crat Una fflated Republcan research

Sample Without Phone Younger

Less Extended Information on Sample Without Phone

Examining Bias in Voter File Between Collection Modes: Overall Research

Problem:

Phone interviewing of voters is getting expensive and has limitations.

- Decreased Response Rate
- Increased Cell Phone Cost
- Hard to show multi-media (Ad-testing)

Online Voter File Match offers solutions:

- Extensive frame information
- Efficiently target registered voter population
- Improved participation rates
- Lower Cost

How should we start doing surveys on the voter file in both modes?

Research:

Used National Data Trust file N=500 Online Matched

- N=500 CATI Non-Matched
- 300 Landline, 200 Cellphone Topic was Tax Reform

Match Criteria

Full Name, Address (Zip+4), Year of Birth* (I) Full Name, Address (Zip+4), Year of Birth (S) Full Name, Zip Code, Year of Birth (M) Household Matches Last Name, Address (Zip+4)

Found important differences in frames

Found important differences in questionnaire modes



Examining Bias in Voter File Between Collection Modes: Online Panel Low Coverage vs. Phone Low Response



Examining Bias in Voter File Between Collection Modes: Need to Quota on Age and Gender

No Partisan Bias in Frame



Online Frame Slightly More Likely to Vote

All R e co rd s	.% 20%	1	5%	12	2%	7%	6% 7	/% 10	% 14%	6 49
Matchedto Phone Number	1% 20%	1	6%	12	2%	7% (5% 8	% 109	% 15%	6 5%
Matched to Online Panel	1%14%	15%	13	%	8%	7%	8%	11%	16%	6%
0%-10 %	10 %	20%	20	%3	0%		30 %4	0%	■ 40 %5	5 0%
50 %6 0%	60 %	70%	70	%8	0%		80 %9	0%	■ 90 %	L 00%

researcn



Online Frame Younger and More Female

Examining Bias in Voter File Between Collection Modes: Mixture of Modes Gets Better Income Representivity

10 0% More Likely 5% 4% to Respond 90% 80% 14% 15% 70 % 12% 16% 60% 24% 17% 50 % 19% More Likely 40 % to Join 27% 22% 30 % 16% 20% 10 % 21% 20% 12% 0% Census Phme* Online ■ \$2 5,0 00 to less th an \$ 50 0 00 ■ Le ss th an \$25,000 ■ \$5 0,0 00 to less th an \$ 75 0 00 ■ \$1 25.000 to b ss than \$1 50,000 ■ \$7 5,0 00 to less th an \$ 10 0,00 0 ■ \$1 00,000 to b ss than \$1 25,0 00



* 14% Item Non-Response

Examining Bias in Voter File Between Collection Modes: Can be Important Mode Effects in the Questions

No Stated Middle Ground

Would you say that the new federal tax law is, in general, a good thing or a bad thing for the country?*

Explicitly Stated Middle Ground

Do you believe that the new federal tax law will increase your taxes, reduce your taxes, or will it not have much of an impact either way?**



* Answer options good/bad rotated in question/answer options. However "Unsure" not read in question, only seen in answer options.





Examining the Accuracy of Likely Voter Models: Overall Research

Problem:

Elections can be great influenced by who turns out to vote. Lots of different models out there:

- Probabilistic Models or Cutoff Models
- Based on Stated Intention or Past Behavior

They can all lead to different results



Which likely voter model is more accurate?



Research:

Used Aristotle Voter File in Virginia and Alabama

Dual Frame RDD/Voter File Survey

- Three surveys (combined N=1,971)
- Tried to match RDD sample back to voter file on first/last name, zip, age.

Compared 9 Methods: Cutoff Models:

- Self-Predict (Low-Cutoff)
- Self-Predict (High-Cutoff)
- Classification Tree Cutoff
- Habitual or highly engaged
- 2016+Probably/Certain
- 2/3 Elections

Cutoff Models:

- Self-Predict (Weighted)
- Classification Tree (Weighted)
- Past-Vote Weighted Model

Examining the Accuracy of Likely Voter Models: All Models Had Similar Individual Accuracy

Voter Turnout Estimate



Individual Level Accuracy of Likely Voters





Examining the Accuracy of Likely Voter Models: Self-Stated Models Had Better Aggregate Estimates

Candidate Prediction Bias

Difference in Vote Choice of Likely Voter Electorate Compared to Baseline of Validated Voters



Partisan Bias in Modeled Turnout





Relational validity: Overall Research

Problem:

A lot of new sampling techniques that don't rely on probability rules (non-probability samples).

We want to find a way to measure the ability to capture true bi-variate relationships (especially non-linear ones).

Earth Mover Distance (EMD) offers solutions:

- Allows for non-linear relationships
- Measures the amount of 'work' to bring back to known relationship
- Identifies the most efficient 'flow' to remove bias

Can we use Earth Mover Distance as a proxy for bias measurement in non-linear bi-variate relationships?

Research:

1) Simulated skewed distributions between age and party support.

25 samples of 1000 voters for each source

2) Applied technique to relationship between first election they voted in and partisanship

- Surveys from 2004-2017
- White registered voters
- 42 surveys
- 17 sample vendors
- 846,024 respondents

Able to rank vendors better on ability to capture the bi-variate non-linear relationship between first election voted and partisanship.



Relational validity: Simulations show skewed sampling -> higher EMD

- Voting propensity tied to Age
 - Increases the EMD score between age and partisanship when either variable moves

	Over Sample Younger	No Age Skew in Sampling	Over Sample Older
Over Sample Active Voters	6.49	6.20	6.21
No Activity Skew in Sampling	8.00	1.68	7.95
Under Sample Active Voters	7.76	7.73	7.07

EMD Lower When No Sampling Skew



Relational validity: Clear Differences by Vendor in Capturing Benchmark

Aggregate Probability Sampling Benchmark

Relationship by Vendor





Relational validity: Final Vendor and Technique Rankings

Project	EMD (centered)	Diff. from Top Rank	EMD rank (centered)
Survey Monkey - Survey Monkey Front Page	0.664	0.000	1
Anonymous - Facebook	1.268	0.604	2
CCES - YouGov	1.491	0.827	3
NAES - Knowledge Networks	1.593	0.929	4
Pewl - American Trends Panel	1.752	1.088	5
ANES - Mail push to web	2.151	1.487	6
ANES - Knowledge Networks	2.459	1.795	7
Miscellaneous - SSI Panel	2.659	1.995	8
CCAP - YouGov	2.736	2.072	9
Broockman et al - Mail push to web	2.765	2.101	10
Lucid - Lucid	3.269	2.605	11
TAPS - Knowledge Networks	4.788	4.124	12
Voter Study Group - YouGov	4.892	4.228	13



Verifying Voter Registration Records: Overall Research

Problem:

Many people take voter file data as truth, but there is error in administration records that is normally ignored.

Here are some of the sources of potential error:

- Concept validity
- Measurement error
- Processing error
- Temporal error

Research:

Used Florida state record file with L2 appends

Sampled 60,000 numbers from FL February, 2017 Voter file

L2 appended phone numbers (59.8%)

Conducted Telephone interviews (N=401)

- 6.5% Response Rate
- Asked Voter Registration question
- If different than what was on the Voter File
 - Got correct information
 - Asked reason for the discrepancy

Substantial error rate (17%) in at least one field on file.

How much error is expected in voter file records?



Verifying Voter Registration Records: Small, but Significant Errors in Voter File



There in non-negligible error in voter records.

Many voter file vendors will clean these up.

- Aristotle
- Catalist
- Data Trust
- i360
- L2
- Target Smart

Varies quite a bit from state to state.

Voter Registration forms need to be updated for easier use (treat like a survey).

Online Panels Matching to Voter Files: Overall Research

Problem:

Many surveys moving to online including political for cost savings and increased capabilities.

However, online panels tend to have low coverage of the overall voter file. This can lead not only to potential

coverage bias but low feasibility. To improve feasibility some have allowed household matches rather than just individual matches to be used in the online research.

Research:

Used National Aristotle file

Compared the following sample frames

- Full Voter File
- Voter File with Phone Number
- Online Panel Individually Matched Full Name, Address (Zip+4), Year of Birth* (I)
 Full Name, Address (Zip+4), Year of Birth (S)
 Full Name, Zip Code, Year of Birth (M)
 Last Name, Address (Zip+4)
- Online Panel Household Matched Last Name, Address (Zip+4)

How does the full voter file compare to either individual or household matches in an online panel sample?

Don't need quotas on partisanship, but maybe on Vote Propensity.



Online Panels Matching to Voter Files: Political Ideology not Skewed in Online Panel

Democrat Affinity Score Comparison

Online PanelIn dividualOn Iy	10%	10%	9%	9%	9% 1	L 0% (10%	10%	11%	12%
Online Panel Tota I	12%	11%	10%	9%	9%	9%	9%	10%	10%	11%
MatchedtoPhoneNumber	12%	11%	5 10%	6 9%	9%	9%	9%	9%	10%	11%
		_								
Ful File	11%	10%	10%	10%	10%	10%	10%	10%	10%	10%

Republican Affinity Score Comparison

Online PanelIn dividualOn Iy	11%	10%	10%	10%	9%	9%	9%	10%	11%	11%
Online Panel Tota I	10%	9%	9%	9%	9%	9%	10%	10%	12%	12%
Matched to Phone Number	11%	9%	9%	9%	9%	9%	10%	11%	12%	13%
Ful File	10%	10%	10%	10%	10%	5 10%	6 10%	5 10%	11%	11%

Party Code Distribution in Matched Frame





Online Panels Matching to Voter Files: Online Panels a Little More Engaged



2020 General Election





2020 Primary Election



21

0 h O

Final Thoughts/Considerations

Remember to respect privacy when using Voter Files

Data will get more and more integrated (Ad Tracking, Geofencing, ect.)

Passive/Big Data from Voter Files likely to supplement, not replace surveys techniques Continue to watch for skews in coverage and response rate for potential bias Never going to be perfect, will continue to improve

