



# That's a Long Survey!

## Using Split-Questionnaire Design to Reduce Respondent Burden in a State Health Survey

Cameron McPhee | SSRS

Cordelia Horch | SSRS

YuChing Yang | University of California – Los Angeles

Jiangzhou Fu | University of California – Los Angeles



# Acknowledgements/Disclaimer

We thank the **University of California – Los Angeles Center for Health Policy Research** for allowing the use of California Health Interview Survey data for this presentation.

Views expressed in this presentation are those of the presenter and not the study sponsor nor any organization with which the authors are affiliated.





# Background

# Surveys are Long...Especially State and Federal Surveys



1. Length Impacts Burden
2. Burden Impacts Cooperation
3. Cooperation Impacts Data Quality and Cost



california  
health  
interview  
survey

> 20,000 unique  
respondent per year  
>700 Unique Survey  
Questions

### Average Adult Survey Length in CHIS:

- 2021 was 49 Minutes
- Web Mode Average is 47 Minutes
- Phone Mode Average is 69 Minutes

Non-English Surveys  
Took Longer  
(~58 Minutes)

Individual Section  
Timings Ranged from  
Less than 1 Minute to  
About 6 Minutes

# Split Questionnaire Designs (*a.k.a. Matrix-Sampling*)

## Methods

- Respondents (rows) and items (columns) are *both* “sampled” from a conceptual complete population data matrix
- Ideally creates missing at random (MAR) or even MCAR data distributions

## Benefits

- Recover univariate and multivariate distributions with, often, minimal loss in precision (Raghunathan & Grizzle, 1995), though findings are mixed (Axenfeld, et., al., 2022)
- Reduce measurement error associated with longer questionnaires (Peytchev & Peytcheva, 2017)

## Constraints


- Must include good predictors of the split items in core modules (Thomas, et. al., 2006)
- Assignment of items to modules can impact the bias introduced through imputation (Axenfeld, et., al., 2022)



# Split-Questionnaire Design (SQD) and CHIS

# Research Questions

## Is there a modular (split-questionnaire) structure that could be implemented for the CHIS that...

- Meaningfully **reduces burden** on respondents?
  - Preserves the **logical flow** of the survey instrument
  - Includes “core” questionnaire items that are **highly correlated** to the items being imputed
  - Maintains all **weighting** items in the core questionnaire
  - Produces unbiased estimates overall and for key subgroups relative to the full questionnaire design
- 
- **Allows for the construction of a usable, fully-imputed, publicly releasable data file (i.e., does not require multiple imputation)**
  - Will not attenuate variance estimates relative to the full questionnaire design



# Methods: Form Design

- 1 Created Modules:** Divided the CHIS questionnaire into topic clusters (modules) of approximately similar numbers of items
- 2 Identified Core Items:** For each topic cluster, identified “core” items
  - Often these were “routing” items that determines the follow-up questions a respondent receives
- 3 Built Forms:** Created 28 unique “forms”, each included
  - All demographic survey questions
  - All items identified as core items for all modules
  - All non-core items for two modules
- 4 Selected Critical Variables:** Chose 21 “key” survey items for simulation and analysis



# Modules, Routing Items, and Dependent Items

	1. Health Conditions and Disability	2. Smoking, Alcohol, and Drug Use	3. Health and Sexual Behavior
<b>CORE ITEMS</b>	<ul style="list-style-type: none"> <li>• General health</li> <li>• Diagnoses: Asthma, Diabetes, etc.</li> <li>• COVID related questions</li> <li>• Tested for colon cancer</li> <li>• Vision/hearing difficulties</li> </ul>	<ul style="list-style-type: none"> <li>• 100+ Cigarettes</li> <li>• E-cig. or vape use</li> <li>• Chewing tobacco/Snuff (30 days)</li> <li>• Marijuana/CBD use</li> <li>• Heroin (12 months)</li> <li>• Alcohol (ever)</li> </ul>	<ul style="list-style-type: none"> <li>• How often eat fruit</li> <li>• Importance of genetics and medical care</li> <li>• # of firearms in home</li> <li>• # of sexual partners (12 months)</li> <li>• Sexual orientation</li> <li>• Ever used PrEP</li> <li>• Ever tested for HIV</li> <li>• Ever received HPV vaccine</li> </ul>
<b>DEPENDENT ITEMS</b>	<ul style="list-style-type: none"> <li>• COVID vaccine receipt</li> </ul>	<ul style="list-style-type: none"> <li>• Person around you smokes/vapes</li> <li>• Days drank 4+ alcoholic drinks</li> </ul>	<ul style="list-style-type: none"> <li>• # sweet beverages (past month)</li> <li>• Importance of environmental &amp; behavioral factors to health</li> <li>• Was offered HIV test</li> </ul>

# Modules, Routing Items, and Dependent Items

	4. Health Insurance Coverage	5. Health Insurance Detail	6. Health Care
<b>CORE ITEMS</b>	<ul style="list-style-type: none"> <li>• Medicare/Medi-CAL</li> <li>• Employer Insurance</li> <li>• Private insurance</li> <li>• CHAMPUS/CHAMP-VA, or military insurance</li> <li>• Other government health insurance program</li> <li>• Other health insurance</li> </ul>	<ul style="list-style-type: none"> <li>• HMO</li> <li>• High deductible</li> <li>• Continual insurance (12 months)</li> <li>• Reason for uninsurance</li> <li>• Reached plan limit</li> <li>• Reason not enrolled in Medi-CAL</li> </ul>	<ul style="list-style-type: none"> <li>• Usual place for care</li> <li>• ER use</li> <li>• Hospital stays</li> <li>• # of doctor's visits</li> <li>• Telehealth use</li> <li>• Difficulty understanding physician</li> <li>• Difficulty/delay receiving medication</li> <li>• Pregnancy status/plans</li> <li>• Mammogram</li> <li>• Dental visits</li> <li>• Racial barriers to care</li> <li>• Need for mental health care</li> </ul>
<b>DEPENDENT ITEMS</b>	<ul style="list-style-type: none"> <li>• Monthly cost of health plan</li> </ul>	<ul style="list-style-type: none"> <li>• Prescription drug coverage</li> <li>• Deductible over \$2,000</li> <li>• Previously had health coverage</li> </ul>	<ul style="list-style-type: none"> <li>• Time since last checkup</li> <li>• Telehealth in past 12 months</li> <li>• Dental insurance</li> </ul>

# Modules, Routing Items, and Dependent Items

	7. Psychological Distress/Mental Health	8. Employment, Housing, & Earnings
<b>CORE ITEMS</b>	<ul style="list-style-type: none"> <li>• Felt nervous, hopeless, restless, depressed, everything was an effort, worthless (past 30 days)</li> <li>• Experienced hazardous climate event</li> <li>• Intimate partner violence</li> <li>• Live with anyone depressed or mentally ill</li> <li>• Able to talk about feelings growing up</li> <li>• # times stopped by police (past 3 years)</li> <li>• Suicidal thoughts (self/close friends)</li> </ul>	<ul style="list-style-type: none"> <li>• Work hours</li> <li>• Length at job</li> <li>• Earnings, income, child support, worker's comp, Social Security/pension</li> <li>• Awareness of CA FMLA laws</li> <li>• Taken paid leave for more than 2 weeks</li> <li>• Receiving TANF or CalWORKS</li> <li>• Housing unit type, tenure, length at address</li> <li>• Help neighbors, neighbors get along, neighbors can be trusted</li> <li>• Volunteered in community</li> <li>• Did not apply for services due to self/family immigration status</li> </ul>
<b>DEPENDENT ITEMS</b>	<ul style="list-style-type: none"> <li>• Phone/computer use (per day)</li> <li>• Physical abuse from intimate partner</li> <li>• Importance of providers asking ACEs</li> <li>• 2 non-parent adults involved in childhood</li> <li>• Ever been arrested</li> </ul>	<ul style="list-style-type: none"> <li>• Feelings about current housing situation</li> <li>• Feel safe in neighborhood</li> <li>• Why not vote in most recent election</li> </ul>

# Form Splits

## Form 1

- Demographics
- Routing Items (all)
- Module 1
- Module 2

## Form 2

- Demographics
- Routing Items (all)
- Module 1
- Module 3

## Form 3

- Demographics
- Routing Items (all)
- Module 1
- Module 4

## Form 4

- Demographics
- Routing Items (all)
- Module 1
- Module 5

## Form 5

- Demographics
- Routing Items (all)
- Module 1
- Module 6

## Form 17

- Demographics
- Routing Items (all)
- Module 3
- Module 7

## Form 18

- Demographics
- Routing Items (all)
- Module 3
- Module 8

## Form 19

- Demographics
- Routing Items (all)
- Module 4
- Module 5

## Form 20

- Demographics
- Routing Items (all)
- Module 4
- Module 6

## Form 21

- Demographics
- Routing Items (all)
- Module 4
- Module 7

## Form 24

- Demographics
- Routing Items (all)
- Module 5
- Module 7

## Form 25

- Demographics
- Routing Items (all)
- Module 5
- Module 8

## Form 26

- Demographics
- Routing Items (all)
- Module 6
- Module 7

## Form 27

- Demographics
- Routing Items (all)
- Module 6
- Module 8

## Form 28

- Demographics
- Routing Items (all)
- Module 7
- Module 8

# Methods: Simulation and Analysis

## Created 50 Replicates of the CHIS:2021 Adult Data File with all Demographic, Routing, and Outcome Variables

- Randomly assigned cases to 1 of 28 forms for each replicate
- Set to missing responses for all questions not included on the assigned form

### Phase 1

1. For each dependent variable, use CART model (R ctree) to predict variable value using demographics and all routing items from all sections *other* than the one for that item
  - Export final “node” from the CART model, and assign to all cases in the full file
2. Run sequential hotdeck imputation on missing cases of each dependent variable using section-specific routing items and final CART node to form imputation classes
3. Repeat imputation (allowing imputed cases to donate values) for remaining missing data, under the same model



# Methods: Simulation and Analysis (continued)

## Phase 2\*:

- Once all dependent variables have been imputed once
  1. Blank imputed values (one variable at a time) from all cases originally missing
  2. Rerun *CART and* hotdeck using all variables (including other imputed variables)
  3. Repeat imputation process 5 times (see Marker, Judkins, and Winglee, 2002).



# Results: Form Length

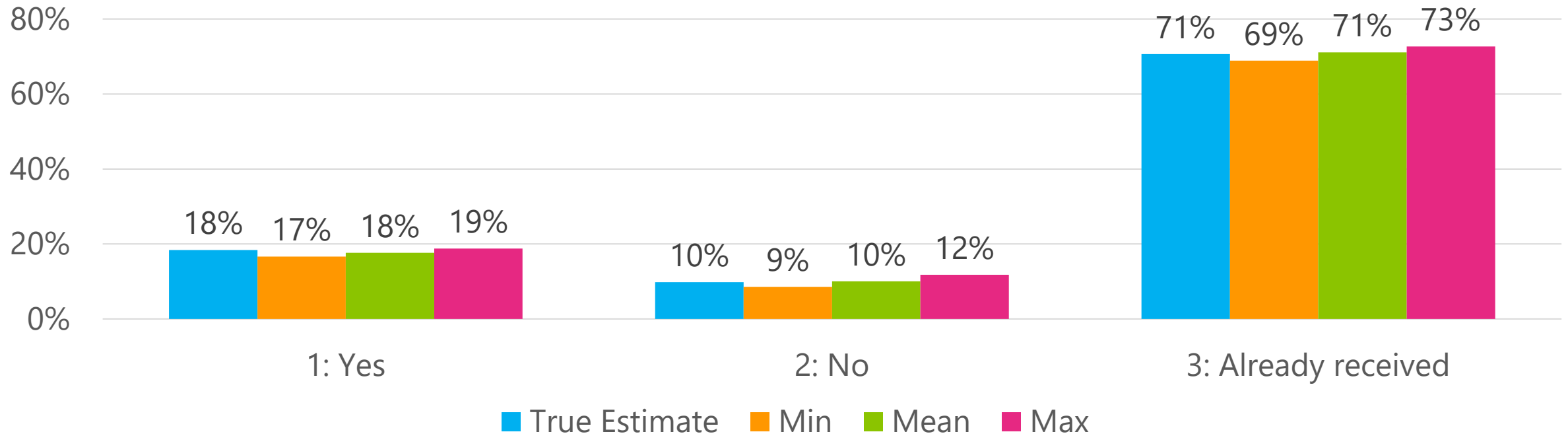
- All forms theoretically would include approx. 80 demographic items & 95 section routing items
- Modules range from 25-89 items
- Form length ranges from 230-345 items
- Respondents complete approximately 10-12 questions per minute (on average)
- New design would yield a survey approximately 25 minutes on average (20-30 minutes depending on the form)
- **Near 50% reduction in survey length!**





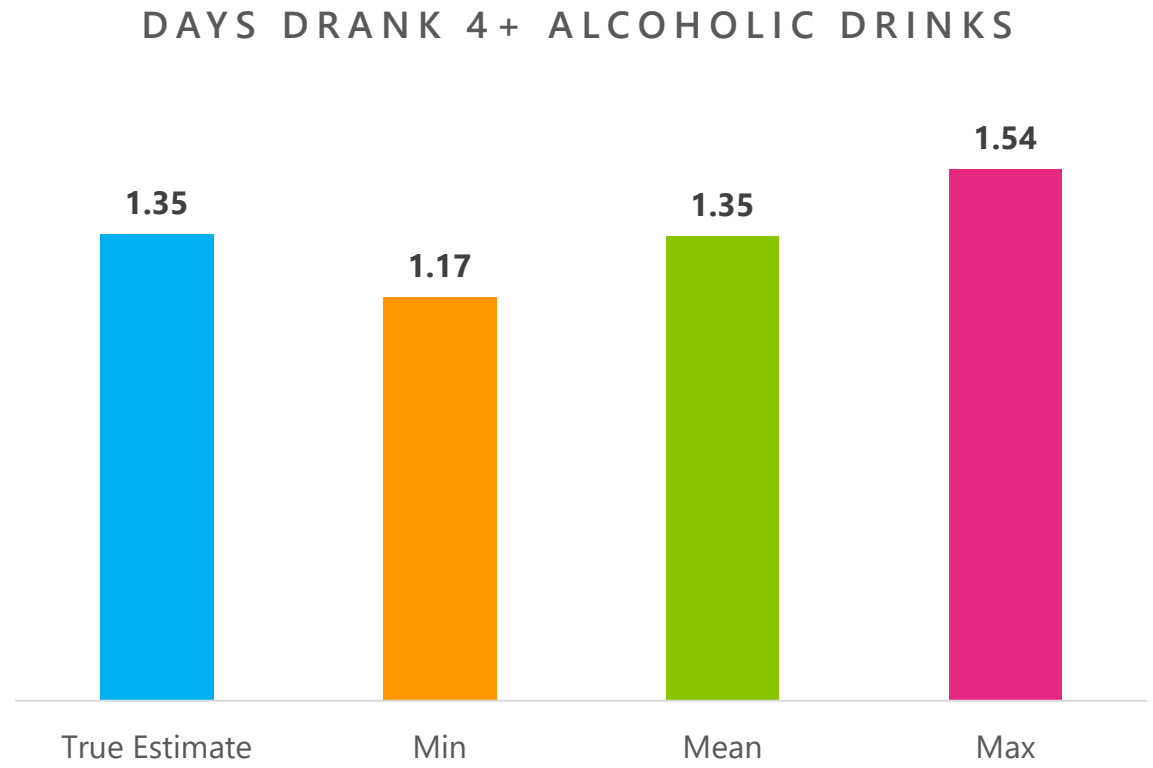
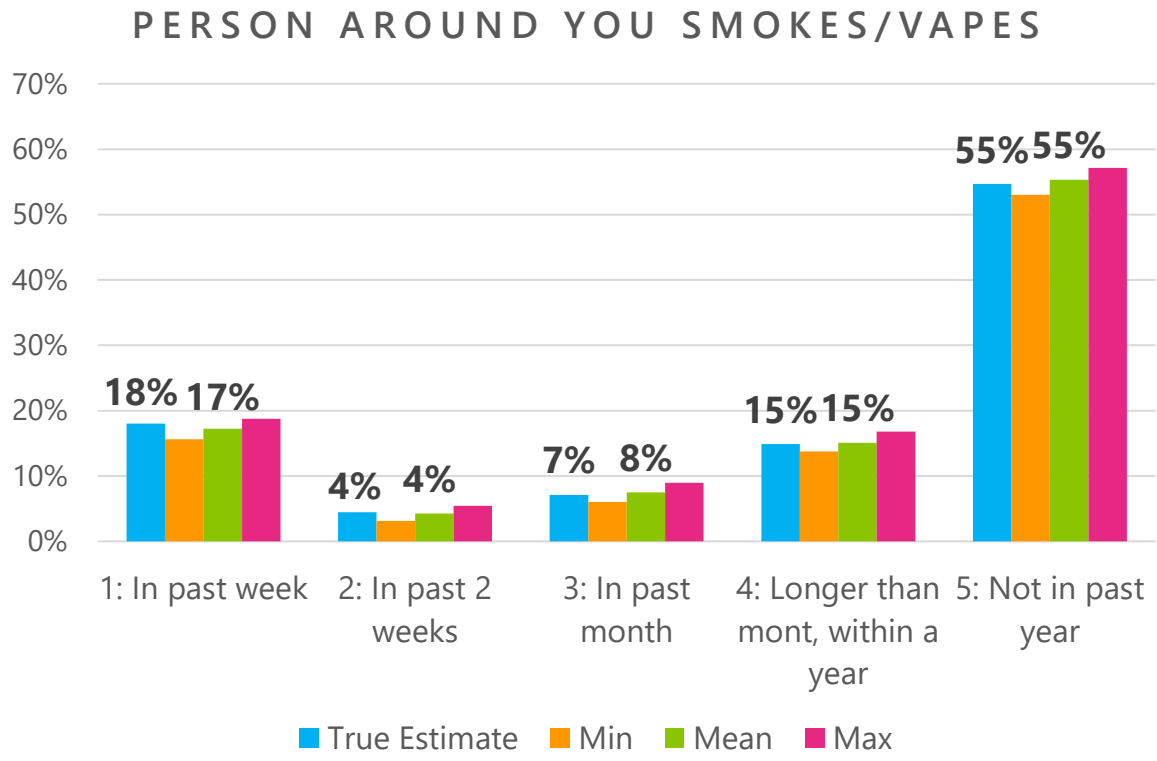
# Health Conditions and Disability Full Sample Estimates

## GOT COVID VACCINE



# Smoking, Alcohol, and Drug Use

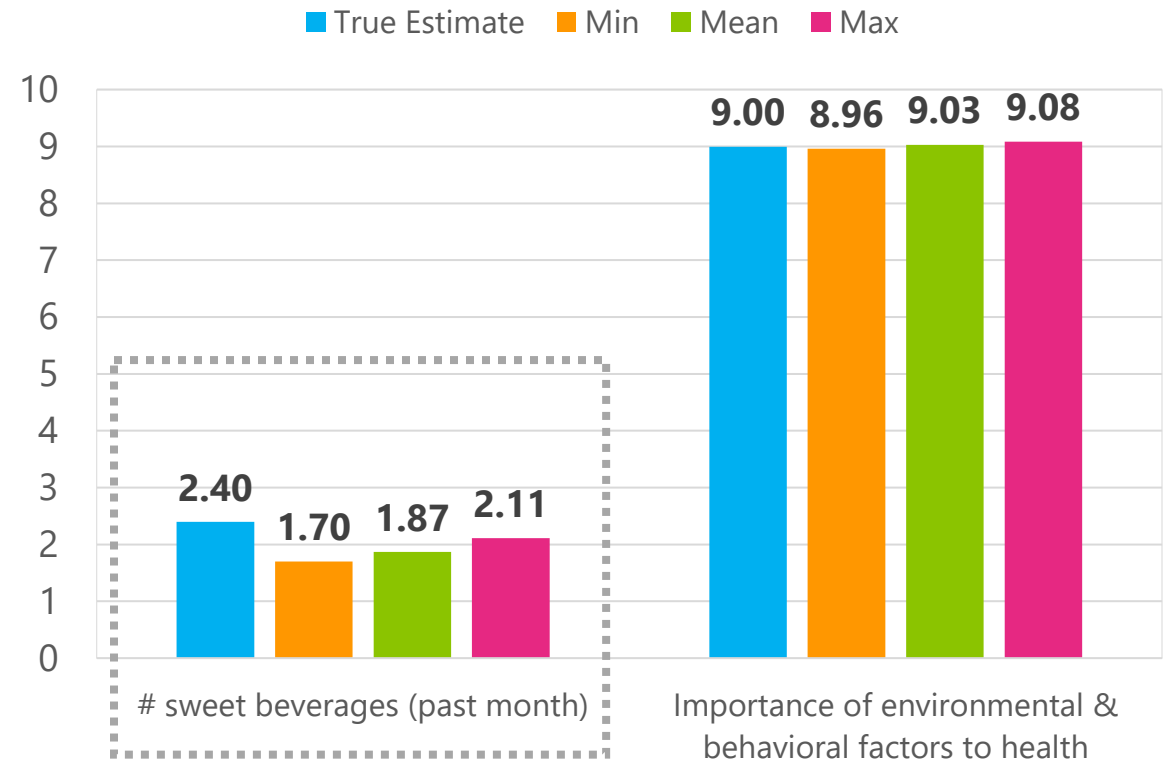
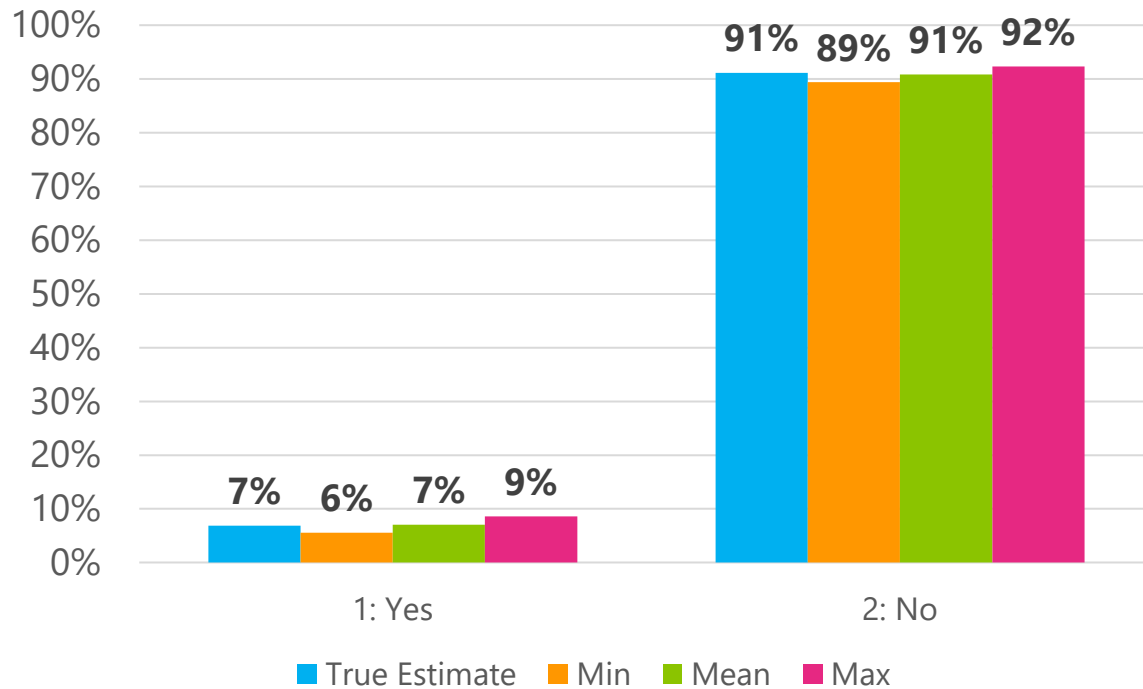
## Full Sample Estimates



# Health and Sexual Behavior

## *Full Sample Estimates*

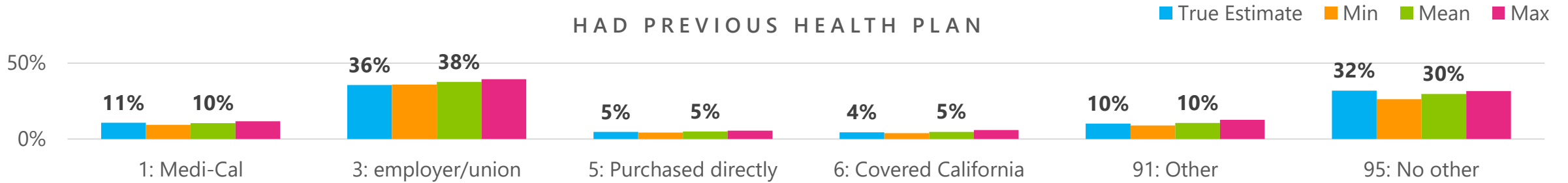
OFFERED HIV TEST



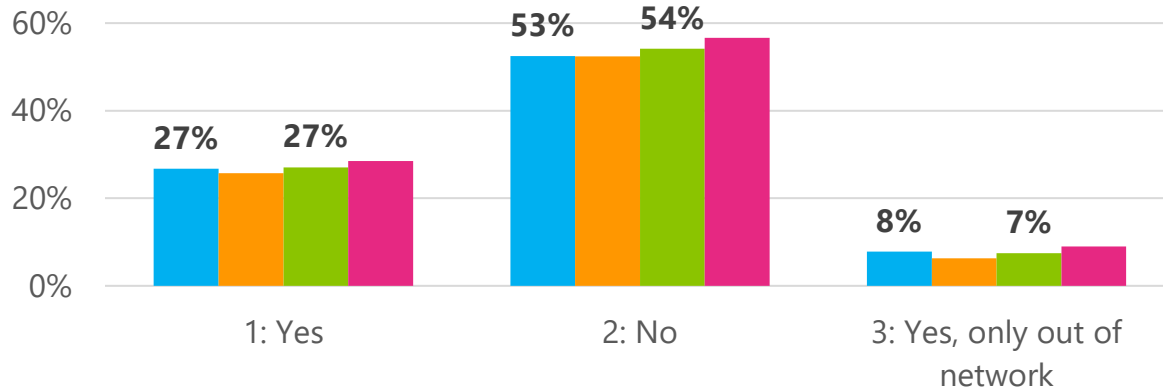
# Health Insurance Coverage & Detail

## *Full Sample Estimates*

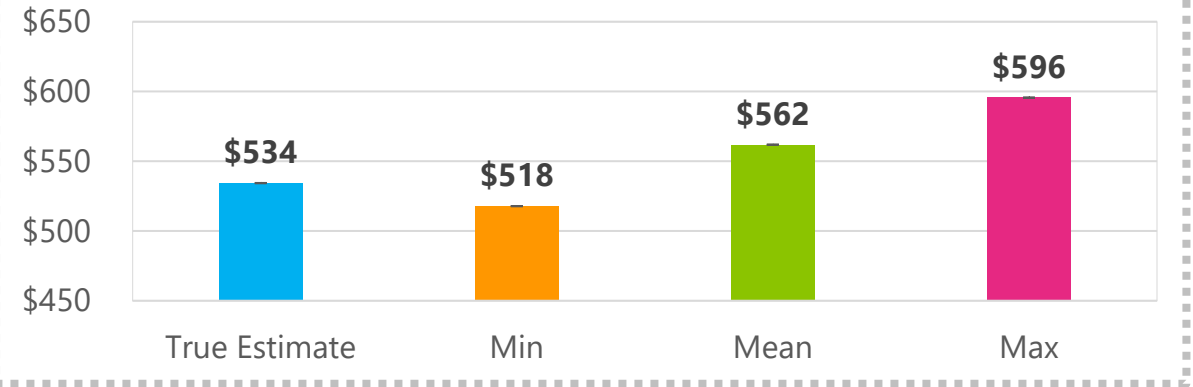
HAD PREVIOUS HEALTH PLAN



DEDUCTIBLE OVER \$2,000

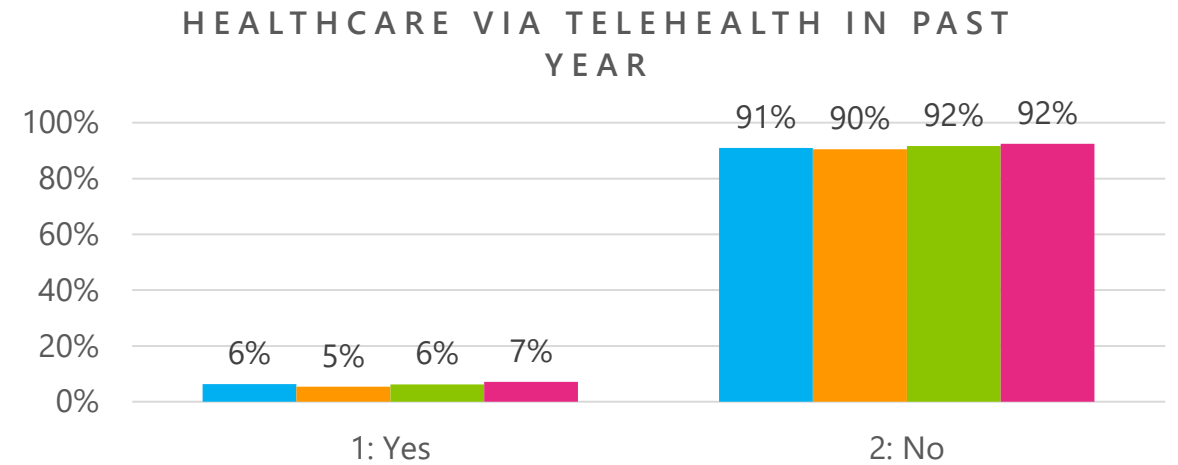
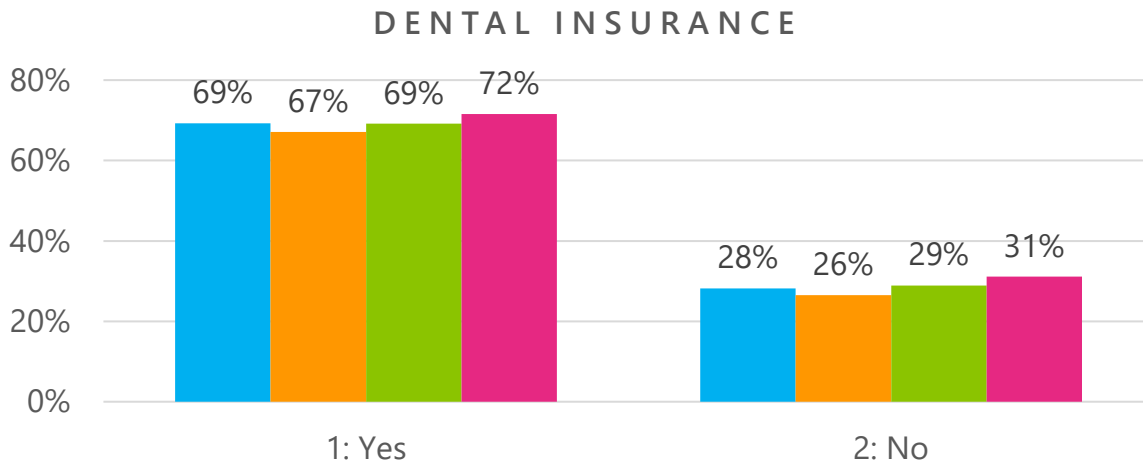
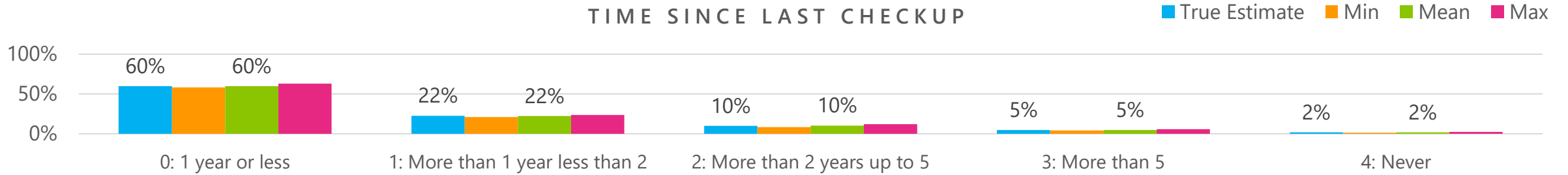


AMOUNT CONTRIBUTE TO HEALTH PLAN PER MONTH



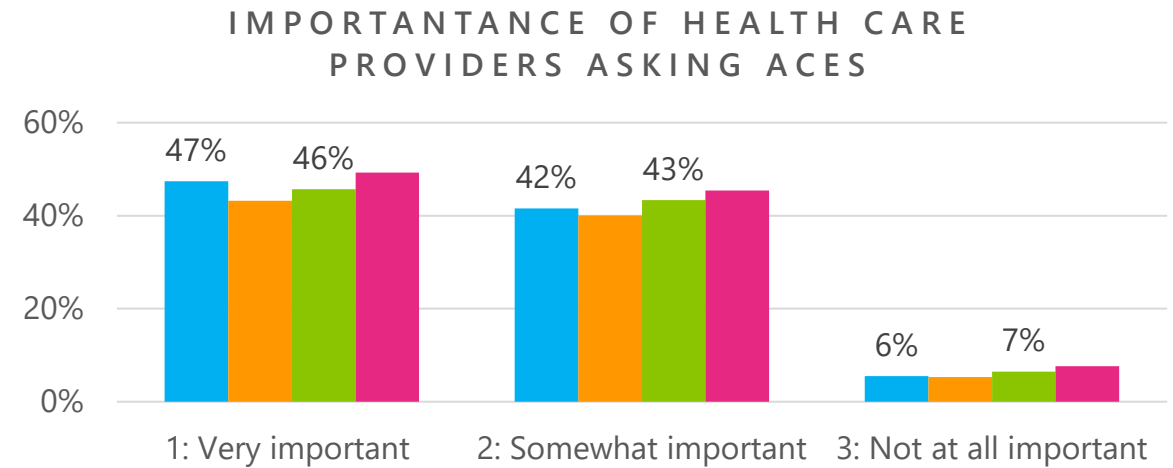
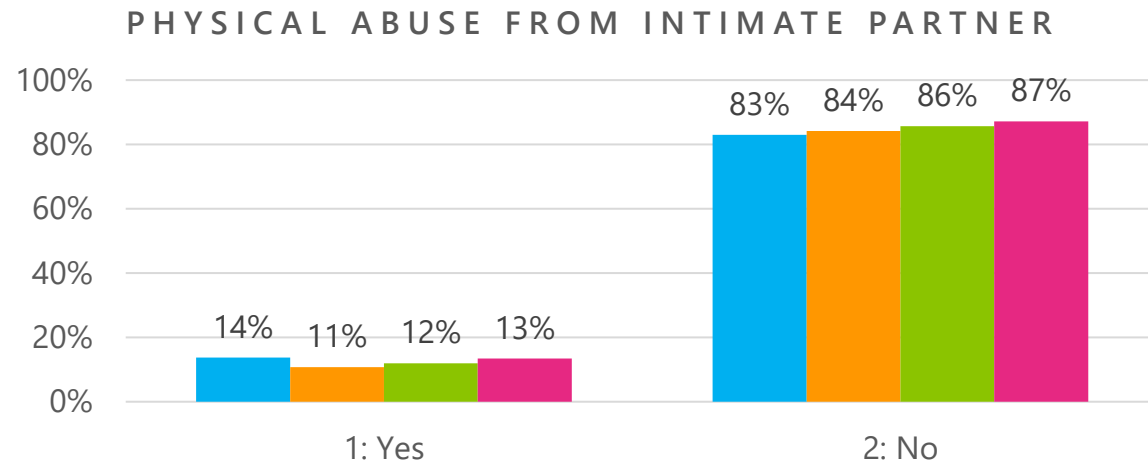
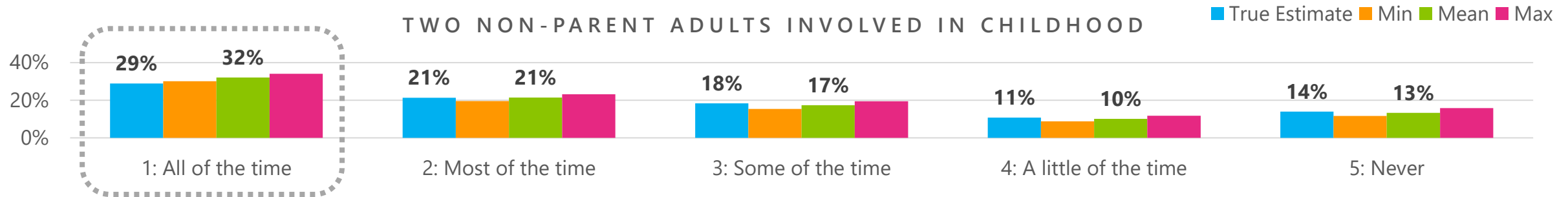
# Health Care

## *Full Sample Estimates*



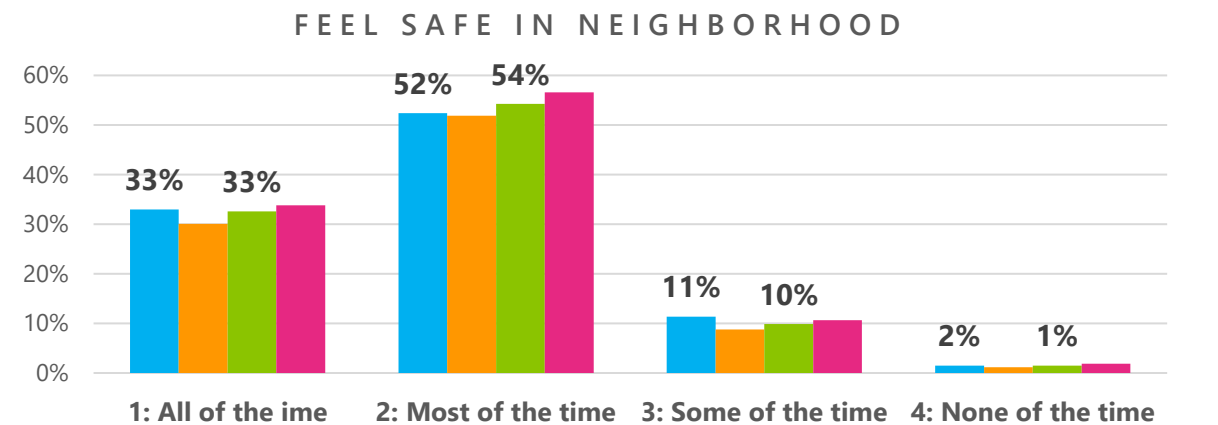
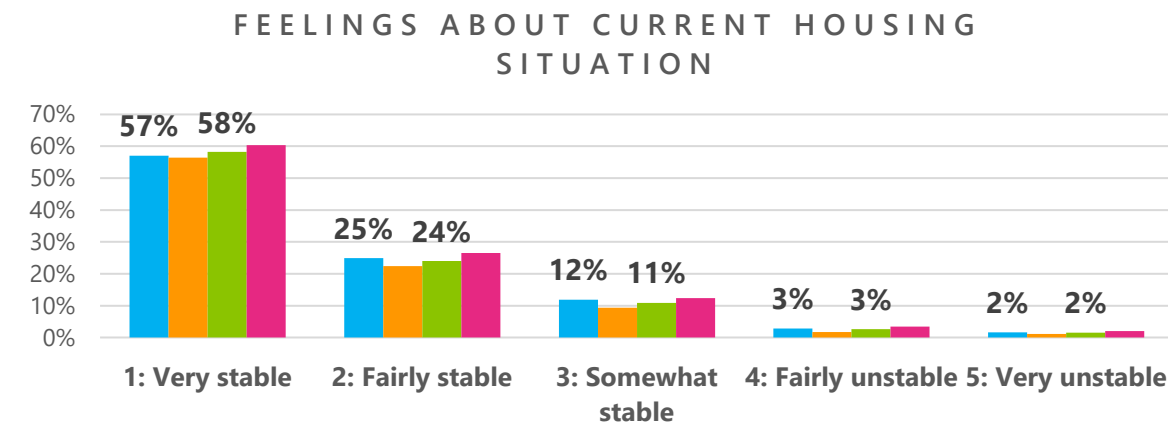
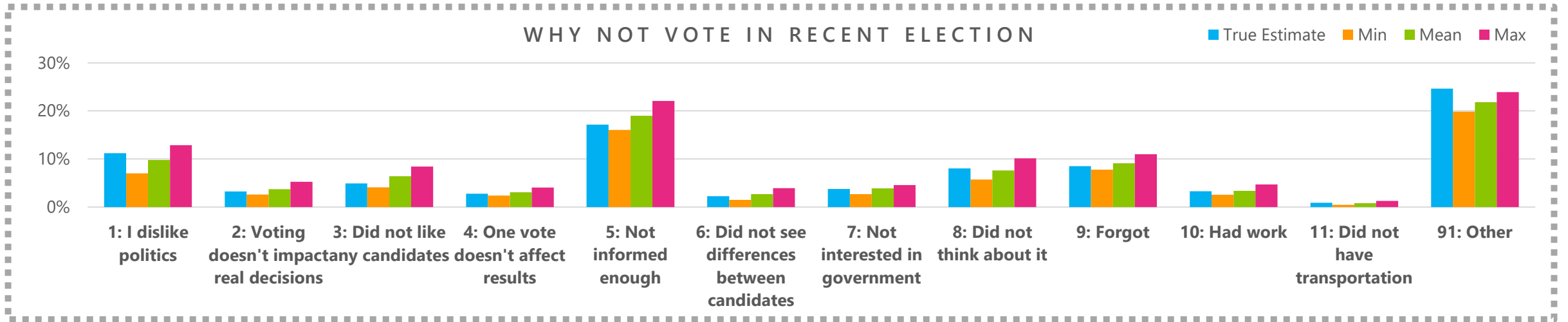
# Psychological Distress/Mental Health

## *Full Sample Estimates*



# Employment, Housing, & Earnings

## Full Sample Estimates



# Variations

	PERCENTAGE ESTIMATES	MEAN ESTIMATES
Average True SE/CV	0.304%	2.23%
Average Imputed SE/CV	0.304%	1.94%
Number of SE's smaller than true value	2,226 (53%)	63 (32%)
Number of SE's greater than true value	1,974 (47%)	133 (68%)



# Summary

## It's HARD!

- Data processing is incredibly time consuming and requires a huge amount of precision coding.

## For the outcomes analyzed here, the imputation led to minimal bias on full sample estimates.

- Items with more categories and continuous variables tend to show more likelihood of bias in imputation.

## The use of hotdeck imputation did not seem to attenuate variances.

- For continuous variables, the SE on the means seemed to be generally larger for the imputed data (which is not a bad thing).
- This may not hold true for subgroup estimates.



# Next Steps

## Extend Findings to Additional Outcomes

Overall and to subgroup estimates

## Form Reduction

- Are all form pairs necessary?
- Can forms be combined such that the longest and shortest modules are paired to narrow range in timings?

## Multivariate Extension

Does the imputation add bias to estimates of relationships?

## Add Cyclic Imputation To Increase Precision

Evaluate whether this reduces the differences in min/max and improves the imputation





# THANK YOU

Cameron McPhee

SSRS Chief Methodologist

[cmcphee@ssrs.com](mailto:cmcphee@ssrs.com)

