

University of Nebraska - Lincoln

DigitalCommons@University of Nebraska - Lincoln

Worlds of Connections Resources

Sociology, Department of

2023

NebrASKa Voices Survey Methodology Report; Including Missing Data Handling and Creating a Raked Weight Variable Using Iterative Proportional Fitting

Joseph C. Jochman

University of Nebraska-Lincoln, joseph.jochman@gmail.com

Julia McQuillan

University of Nebraska-Lincoln, jmcquillan2@Unl.edu

Grace Kelly

University of Nebraska-Lincoln, gkelly@huskers.unl.edu

Patricia Wonch Hill

University of Nebraska - Lincoln, phill3@unl.edu

Meghan Leadabrand

University of Nebraska - Lincoln, megan.leadabrand@gmail.com

Follow this and additional works at: <https://digitalcommons.unl.edu/wrldconnex>



Part of the [Science and Mathematics Education Commons](#)

Jochman, Joseph C.; McQuillan, Julia; Kelly, Grace; Wonch Hill, Patricia; and Leadabrand, Meghan, "NebrASKa Voices Survey Methodology Report; Including Missing Data Handling and Creating a Raked Weight Variable Using Iterative Proportional Fitting" (2023). *Worlds of Connections Resources*. 7. <https://digitalcommons.unl.edu/wrldconnex/7>

This Article is brought to you for free and open access by the Sociology, Department of at DigitalCommons@University of Nebraska - Lincoln. It has been accepted for inclusion in Worlds of Connections Resources by an authorized administrator of DigitalCommons@University of Nebraska - Lincoln.

Nebraska Annual Social Indicators Survey
Two-Wave NebrASKa Voices Surveys
Wave 1: *Summer 2018, Summer 2019, Winter 2019, Winter 2020*
Wave 2: *Summer 2020*
NebrASKa Voices Survey Methodology Report
Including Missing Data Handling and Creating a *Raked Weight Variable* Using Iterative
Proportional Fitting
Prepared January 2021–March 2023
Joseph Jochman, Julia McQuillan, Grace Kelly, Trish Wonch Hill, Meghan Leadabrand

To gain a comprehensive picture of the surveys that this methodology report references, see the Bureau of Sociological Research methodology reports for each survey available at this [link](#). In the original survey reports the Bureau of Sociological Research staff provide descriptions of data collection, sampling and questionnaire design, response rate, data processing, and preliminary data cleaning.

Introduction

This report outlines the process of creating a raked weight variable for the NebrASKa Voices 2020 survey, a second “Wave 2” survey following four different “Wave 1” surveys. The 2020 NebrASKa Voices sample Wave 1 prior Nebraska Annual Social Indicators Surveys (NASIS) come from the NASIS conducted in Summer 2018 (N=116), Summer 2019 (N=162), Winter 2019 (N=55), and Winter 2020 (N=171). Participants in the four Wave 1 surveys were given the option to opt into future research; if they chose to opt in, they became part of the sampling frame for the Wave 2 NebrASKa Voices survey.

The combination of four Wave 1 and one Wave 2 samples complicated the creation of an accurate weight variable. This weighting was necessary due to several factors. First, the initial selection into the sample involved people volunteering to be in the sample from a random sample, thus creating a new sample that may or may not proportionally represent the demographic characteristics of the state of Nebraska (NE). Some groups, for example those who are older, female/women, or white, were more likely to opt in and to complete the survey. Differential propensities to volunteer and to complete the surveys were amplified by the COVID-19 pandemic. To make the two-wave data represent the population, we calculated to account for the selection effects at several time points. Using prior NASIS probability weights for each respondent, we use iterative proportional fitting raking (ipfraking) (Kolenikov 2014, 2019). Ipfraking creates a raked weight variable representative of the American Community Survey (ACS) 2019 NE control totals (i.e., age, female, nonwhite, education). The weight variable allows us to estimate values that represent the Nebraskan population even with the challenges of data collection during the pandemic and multiple Wave 1 samples. Probability weight values (i.e., P_{wate}) were calculated for each respondent based upon their inclusion in their prior NASIS. Note that 29 respondents had missing probability weight values in the final Voices sample. The data for respondents with missing probability weight values could not be imputed, and for accuracy, they were excluded from the sample.

This report is organized into eight (8) steps with associated output:

1. Downloading Voices data and setting initial directories
2. Cleaning/recoding Voices control variables (i.e., age, female, nonwhite, education)
3. Hot deck imputation of Voices control variables (i.e., age, female, nonwhite, education) that were pulled in from previous surveys and imputed for complete cases, except for 29 cases missing data.
4. Selecting/download NE ACS control totals for comparison (i.e., age, female, nonwhite, education)
5. Recoding/checking imputed Voices variable values to match NE ACS control values

6. Setting control matrices using NE ACS totals (N=15,313)
7. Using the ipfraking command (Kolenikov 2014, 2019) to adjust sampling weights
8. Comparing unweighted/weighted values for age, female, nonwhite, education, and other variables

Step 1: Downloading dataset and directories

```
using == "Nebraska Voices 2020 Data_FINAL_NASIS weights added.dta"
.do file == "NE_voices_2021_creating_raked_weights_hotdeck_1.do"
```

Step 2: Cleaning/recoding age, female, nonwhite, and education (Voices/NASIS)

```
.do file == "NE_voices_2021_creating_raked_weights_hotdeck_1.do"
```

a) age (lines 21-57):

Variable	Obs	Mean	Std. Dev.	Min	Max
age_new_V	457	57.25821	15.63492	24	93

b) female (lines 58-72):

Variable	Obs	Mean	Std. Dev.	Min	Max
female_V	454	.592511	.4919092	0	1

c) nonwhite (lines 73-101):

Variable	Obs	Mean	Std. Dev.	Min	Max
nonwhite_V	454	.0682819	.2525073	0	1

d) education (lines 102-160):

d1. Education values from responses in Wave 1 surveys (NASIS 2018-2020)

What is the highest degree you have attained?	Freq.	Percent	Cum.
1. No diploma	2	0.40	0.40
2. High school diploma/GED	37	7.34	7.74
3. Some college but no degree	95	18.85	26.59
4. Technical/associate/junior college	54	10.71	37.30
5. Bachelor's degree	151	29.96	67.26
6. Graduate degree	118	23.41	90.67
.	47	9.33	100.00
Total	504	100.00	

d2. education (Voices 2020)

What is the highest degree you have attained?	Freq.	Percent	Cum.
1. No diploma	5	0.99	0.99
2. High School Diploma/GED	80	15.87	16.87

3. Technical/Associate/Junior College	63	12.50	29.37
4. Bachelor's degree (4yr., BA, BS, RN)	146	28.97	58.33
5. Graduate Degree (Masters, PhD, Law)	107	21.23	79.56
.	103	20.44	100.00
-----+-----			
Total	504	100.00	

d3. recoding education in 2018- 2020 NASIS variable (i.e., combine H/some college) to match Voices (2020)

RECODE of degr (What is the highest degree you have attained?)	Freq.	Percent	Cum.
-----+-----			
1. No diploma	2	0.40	0.40
2. HS/GED	132	26.19	26.59
3. Associate	54	10.71	37.30
4. BA	151	29.96	67.26
5. Grad	118	23.41	90.67
.	47	9.33	100.00
-----+-----			
Total	504	100.00	

d4. replacing missing values of education in Voices (2020) with NASIS (2018-2020) values if applicable

RECODE of Q46 (What is the highest degree you have attained?)	Freq.	Percent	Cum.
-----+-----			
1. No diploma	4	0.79	0.79
2. HS/GED	118	23.41	24.21
3. Associate	72	14.29	38.49
4. BA	169	33.53	72.02
5. Grad	130	25.79	97.82
.	11	2.18	100.00
-----+-----			
Total	504	100.00	

e) summarizing recoded variables:

Variable	Obs	Mean	Std. Dev.	Min	Max
-----+-----					
age_new_V	457	57.25821	15.63492	24	93
female_V	454	.592511	.4919092	0	1
nonwhite_V	454	.0682819	.2525073	0	1
educ_cat_V	493	3.614604	1.137655	1	5

Step 3: Hot deck imputation for age, female, nonwhite, and education categories (command requires Stata version 14 or higher) [work done using the University of NE soc-analyzer]

Hot deck imputation replaces missing values in a dataset with values from another randomly selected observation within the dataset. The randomly selected observation is determined based on a set of matching criteria to impute the most plausible values for the missing data.
 .do file == "NE_voices_2021_creating_raked_weights_hotdeck_1.do"

a) code (line 186): `hotdeckvar age_new_V female_V nonwhite_V educ_cat_V, suffix("_m")`

b) output:

```
Number of observations without missing values:448
Number of observations with missing values:56
Imputing age_new_V_m
(47 real changes made)
Imputing female_V_m
(50 real changes made)
Imputing nonwhite_V_m
(50 real changes made)
Imputing educ_cat_V_m
(11 real changes made)
```

c) summarizing variables (these frequencies can be compared to variable frequencies from Step 2: a-e)

Variable	Obs	Mean	Std. Dev.	Min	Max
age_new_V_m	504	57.25794	15.54756	24	93
female_V_m	504	.5892857	.4924523	0	1
nonwhite_V_m	504	.0654762	.24761	0	1
educ_cat_V_m	504	3.625	1.13688	1	5

d) saving data old == "Nebraska Voices 2020 Data_FINAL_NASIS weights added v02 hotdeck"
 note: data saved "old" to read in using Stata 13 [Step 3 work above done using soc-analyzer]

Step 4: Selecting/downloading Census Bureau ACS NE population totals 2019 from Integrated Public Use Microdata Systems

.do file == "ACS_setup_data_NE_Voices_v03_ACS_CONTROL_TOTALS.do"

link: <https://usa.ipums.org/usa-action/variables/group>

note: requires registration

select/download indicator variables for age, female, race/ethnicity, and education

NE totals below used as control totals with ipfraking

Running frequency tables of ACS demographic variables to determine appropriate weights for NASIS and Voices sample demographic characteristics, using Census counts. The *Ns* are inflated because the counts include multiple NASIS baseline datasets. However, the proportions are constant.

a1) ACS age (US 2019) (note: omitting age less than 18: N= 640,382)

-> tabulation of ACS_agecat

RECODE of ACS_age (age)	Freq.	Percent	Cum.
18-24	271,825	10.46	10.46

25-34	386,595	14.87	25.33
35-44	377,284	14.52	39.85
45-54	402,509	15.49	55.33
55-64	483,610	18.61	73.94
65+	677,348	26.06	100.00

Total	2,599,171	100.00	

a2) ACS age (NE 2019) (note: omitting age less than 18: N= 4,453)

-> tabulation of ACS_agecat

RECODE of ACS_age (age)	Freq.	Percent	Cum.
18-24	1,628	10.63	10.63
25-34	2,271	14.83	25.46
35-44	2,211	14.44	39.90
45-54	2,107	13.76	53.66
55-64	2,856	18.65	72.31
65+	4,240	27.69	100.00

Total	15,313	100.00	

b1) ACS female (US 2019)

-> tabulation of ACS_female

RECODE of ACS_sex (sex)	Freq.	Percent	Cum.
0. male	1,259,569	48.46	48.46
1. female	1,339,602	51.54	100.00

Total	2,599,171	100.00	

b2) ACS female (NE 2019)

-> tabulation of ACS_female

RECODE of ACS_sex (sex)	Freq.	Percent	Cum.
0. male	7,567	49.42	49.42
1. female	7,746	50.58	100.00

Total	15,313	100.00	

c1) ACS nonwhite (US 2019)

-> tabulation of ACS_nonwhite

RECODE of |
ACS_race |
(race |

[general version])	Freq.	Percent	Cum.
0. white	2,039,728	78.48	78.48
1. nonwhite	559,443	21.52	100.00
Total	2,599,171	100.00	

c2) ACS nonwhite (NE 2019)

-> tabulation of ACS_nonwhite

RECODE of ACS_race (race [general version])	Freq.	Percent	Cum.
0. white	14,036	91.66	91.66
1. nonwhite	1,277	8.34	100.00
Total	15,313	100.00	

d1) ACS education (US 2019)

-> tabulation of ACS_cateduc

RECODE of ACS_educd (educational attainment [detailed version])	Freq.	Percent	Cum.
1. No diploma	274,448	10.56	10.56
2. High School Diploma/GED	1,279,063	49.21	59.77
3. Technial/Associate/Junior College	217,680	8.37	68.14
4. Bachelor's Degree (4yr., BA, BS, RN)	507,242	19.52	87.66
5. Graduate Degree (Masters, PhD, Law)	320,738	12.34	100.00
Total	2,599,171	100.00	

d2) ACS education (NE 2019)

-> tabulation of ACS_cateduc

RECODE of ACS_educd (educational attainment [detailed version])	Freq.	Percent	Cum.
1. No diploma	1,084	7.08	7.08
2. High School Diploma/GED	8,139	53.15	60.23
3. Technial/Associate/Junior College	1,735	11.33	71.56
4. Bachelor's Degree (4yr., BA, BS, RN)	2,919	19.06	90.62
5. Graduate Degree (Masters, PhD, Law)	1,436	9.38	100.00
Total	15,313	100.00	

Step 5: Recoding/checking hot deck imputed Voices variables to match ACS control values
 using == "Nebraska Voices 2020 Data_FINAL_NASIS weights added v02 hotdeck"
 .do file == "NE_Voices_2021_creating_raked_weights_ipfraking_2.do"

-> tabulation of agecat_new_V

RECODE of age_new_V_m	Freq.	Percent	Cum.
1. 18-24	5	0.99	0.99
2. 25-34	42	8.33	9.33
3. 35-44	68	13.49	22.82
4. 45-54	83	16.47	39.29
5. 55-64	113	22.42	61.71
6. 65+	193	38.29	100.00
Total	504	100.00	

-> tabulation of female_V_m

RECODE of sex (Are you:)	Freq.	Percent	Cum.
0. male	207	41.07	41.07
1. female	297	58.93	100.00
Total	504	100.00	

-> tabulation of nonwhite_V_m

nonwhite_V_ m	Freq.	Percent	Cum.
0. white	471	93.45	93.45
1. nonwhite	33	6.55	100.00
Total	504	100.00	

-> tabulation of educ_cat_V_m

RECODE of Q46 (What is the highest degree you have attained?)	Freq.	Percent	Cum.
1. No diploma	4	0.79	0.79
2. HS/GED	119	23.61	24.40
3. Associate	74	14.68	39.09
4. BA	172	34.13	73.21
5. Grad	135	26.79	100.00
Total	504	100.00	

Step 6: Setting control total matrices using NE ACS totals (see Step 4 using NE totals N=15,313)
.do file == "NE_Voices_2021_creating_raked_weights_ipfraking_2.do"


```

///setting up the totals
capture drop _one
generate byte _one = 1
scalar ACS2019_NE_total_pop = 15313

////age matrix
matrix ACS2019_age = (1628, 2271, 2211, 2107, 2856, 4240)
matrix colnames ACS2019_age = 1 2 3 4 5 6
matrix coleq ACS2019_age = _one
matrix rownames ACS2019_age = agecat_new_V
matrix list ACS2019_age, f(%12.0g)

////sex matrix
matrix ACS2019_sex = (7567, 7746)
matrix colnames ACS2019_sex = 0 1
matrix coleq ACS2019_sex = _one
matrix rownames ACS2019_sex = female_V_m
matrix list ACS2019_sex, f(%12.0g)

////race matrix
matrix ACS2019_nonwhite = (14036, 1277)
matrix colnames ACS2019_nonwhite = 0 1
matrix coleq ACS2019_nonwhite = _one
matrix rownames ACS2019_nonwhite = nonwhite_V_m
matrix list ACS2019_nonwhite, f(%12.0g)

////educ matrix
matrix ACS2019_cateduc = (1084, 8139, 1735, 2919, 1436)
matrix colnames ACS2019_cateduc = 1 2 3 4 5
matrix coleq ACS2019_cateduc = _one
matrix rownames ACS2019_cateduc = educ_cat_V_m
matrix list ACS2019_cateduc, f(%12.0g)

```

Step 7: Using ipfraking and saving data

Kolenikov 2014, 2019

The ipfraking command adjusts survey weights, so that the sample distribution (of Voices data) matches the population distribution (ACS 2019 data) of demographic variables. This command estimates the post-stratification weights based on the variables specified in the command—in this case age, sex, race, and education.

```
.do file == "NE_Voices_2021_creating_raked_weights_ipfraking_2.do"
```

note: 29 missing values on Pwate – values omitted from final raked weight

note: drop missing weight values or replace

```
ipfraking [pw=Pwate], ctotal(ACS2019_age ACS2019_sex ACS2019_nonwhite
ACS2019_cateduc) generate(rakedwgt)
```

```

(29 missing values generated)
(29 missing values generated)
(29 missing values generated)
Iteration 1, max rel difference of raked weights = 833.45421
Iteration 2, max rel difference of raked weights = .62201468
Iteration 3, max rel difference of raked weights = .12246612
Iteration 4, max rel difference of raked weights = .02734436

```

```

Iteration 5, max rel difference of raked weights = .00574776
Iteration 6, max rel difference of raked weights = .0011763
Iteration 7, max rel difference of raked weights = .00023886
Iteration 8, max rel difference of raked weights = .00004853
Iteration 9, max rel difference of raked weights = 9.886e-06
Iteration 10, max rel difference of raked weights = 2.019e-06
Iteration 11, max rel difference of raked weights = 4.129e-07
The worst relative discrepancy of 5.8e-08 is observed for educ_cat_V_m == 5
Target value =      1436; achieved value =      1436

```

Summary of the weight changes

	Mean	Std. dev.	Min	Max	CV
Orig weights	1.0278	.80212	.16787	5.4661	.7804
Raked weights	32.238	62.573	1.8269	925.22	1.941
Adjust factor	35.3077		6.0572	1027.9944	

(29 missing values generated)

Notes on descriptive statistics from Kolenikov 2014:

p. 7 (2014) “In practice, I have encountered increases of this coefficient of variation [CV] between 20% and 100% of the relative scale, or between .2 and 1.5 on the absolute scale, for design effects varying between 1 and 2 in the typical public opinion surveys.”

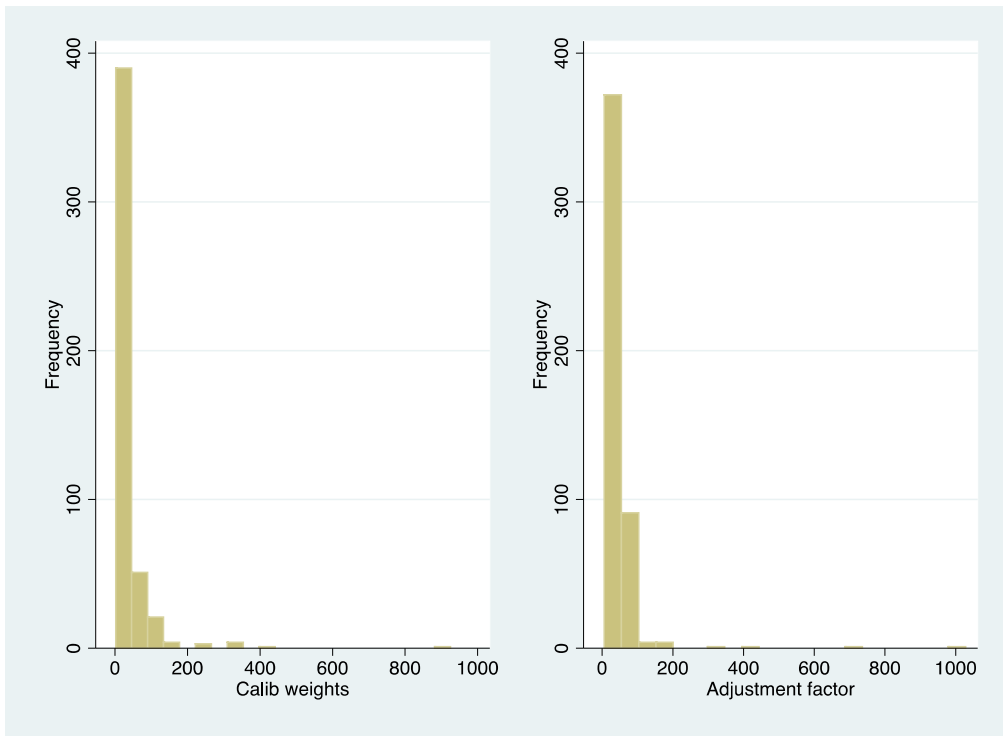
p. 11 (2014) “Besides the internal convergence diagnostics, the weights produced by ipfraking were compared to those produced by survwt and ipfweight as a certification step (Gould 2001), and were found to be identical within numerical accuracy.”

p. 17 (2014) “as expected, the coefficient of variation went up...”

p. 17 (2014) “Generally, we would want to inspect these graphs to see if there [are] any unexpected patterns, such as high outlying values, gaps in the distribution, or concentration near the limits of the weight range.”

Examining the descriptive statistics for the raked weight and the graphic below for the weights for the NebrASKa Voices sample, we conclude that the raked weight is consistent with the expectations described by Kolenikov (2014), who created the weighting software package. In addition, as is evident in the “Step 8” table below, the descriptive statistics with the raked variable for age, gender/female, race and education closely match the ACS proportions. We therefore proceeded with the analyses using the raked weights. Even though the demographic characteristic distributions are distinct between the unweighted and the weighted estimates, the mean values for the network science variables introduced in the NebrASKa Voices survey are similar for the weighted and unweighted estimates.

Associated output graphic:



saving “NE_Voices_2020_data_with_raked_weight_v01.dta”

Step 8: Comparing means for age, female, nonwhite, education, et al.

Variable	Unweighted values				NASIS weight values		Raked weight values		NE ACS totals		US ACS totals	
	M/P ^a	SD	Min	Max	M/P	SD	M/P	SD	M/P	SD	M/P	SD
<i>Age</i>												
Age (continuous)	57.14	15.42	24	91	50.75	15.24	49.77	18.11				
18-24	.01		0	1	.01		.12		.11		.10	
25-34	.09		0	1	.15		.16		.15		.15	
35-44	.13		0	1	.21		.14		.14		.15	
45-54	.16		0	1	.21		.11		.14		.15	
55-64	.23		0	1	.21		.19		.19		.19	
65+	.38		0	1	.22		.27		.28		.26	
<i>Sex</i>												
Female	.59		0	1	.51		.51		.51		.52	
<i>Race/ethnicity</i>												
Nonwhite	.06		0	1	.07		.08		.08		.22	
<i>Education</i>												
HS or less	.23		0	1	.19		.58		.60		.60	
Associate	.14		0	1	.15		.12		.11		.08	
Bachelors	.36		0	1	.41		.20		.19		.20	
Graduate	.27		0	1	.25		.10		.09		.12	
<i>Political affiliation</i>												
Democrat	.35		0	1	.31		.29					
Republican	.41		0	1	.43		.38					
Independent	.21		0	1	.23		.30					
Other party	.03		0	1	.03		.03					
<i>Religious affiliation</i>												
Protestant	.46		0	1	.43		.33					
Catholic	.21		0	1	.20		.20					
Other religion	.16		0	1	.16		.20					
Unaffiliated	.18		0	1	.21		.27					
<i>Network science items</i>												
1a. Heard no	.74		0	1	.75		.73					
1b. Heard yes	.21		0	1	.21		.22					
1c. Heard DK	.05		0	1	.04		.06					
2. Spread disease	2.61	1.10	1	4	2.65	1.11	2.62	1.11				
3. Connections	2.69	1.13	1	4	2.75	1.13	2.68	1.09				
4. Addiction	2.29	1.00	1	4	2.31	1.00	2.37	1.00				
5. Learn more	2.86	.91	1	4	2.85	.91	2.76	.92				
6. Understand health	3.51	.79	1	4	3.56	.74	3.59	.73				
7. Math models	3.54	.79	1	4	3.57	.75	3.45	.97				
8. Improve health	3.00	1.10	1	4	3.01	1.10	3.02	1.10				

N=438

a. Mean/proportion

1. Have you ever heard about network science? (1=yes)
2. How involved do you think network science is in studying the spread of contagious disease? (range 1 “not at all involved” – 4 “very involved”)
3. How useful are network science models for seeing connections that are important for health (e.g., among jobs, food, and schools)? (range 1 “not at all useful” – 4 “very useful”)
4. How much is network science helpful for understanding addiction experiences? (range 1 “not at all helpful” – 4 “very helpful”)

5. How interested are you in learning more about network science? (range 1 “not at all interested” – 4 “very interested”)
 6. How valuable is research on connections among people for understanding human health? (range 1 “not at all valuable” – 4 “very valuable”)
 7. How important are mathematical models of people being near each other (e.g., in schools or workplaces) to understand the spread of contagious disease? (range 1 “not at all important” – 4 “very important”)
 8. How important is network science research for improving public health? (range 1 “not at all important” – 4 “very important”)
- * Items 2-8 DK coded as 1.5

Funding

Research reported in this publication was supported by the National Institute of General Medical Sciences of the National Institutes of Health under award R25GM129836 (“Worlds of Connections: Engaging Youth with Health Research through Network Science and Stories in Augmented Reality”). The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

References

- Andridge, Rebecca R. and Roderick J. A. Little. 2010. A review of hot deck imputation for survey non-response. *International Statistical Review*, 78(1): 40-64.
- Bureau of Sociological Research. 2020. Nebraska Voices 2020 Survey Methodology Report. University of Nebraska-Lincoln.
- Heeringa, Steven G., Brady T. West, and Patricia A. Berglund. 2017. “Applied Survey Data Analysis, Second Edition”. Boca Raton, FL: CRC Press.
- Kolenikov, Stanislav. 2019. Updates to the ipfraking system. *The Stata Journal*, 19(1): 143-184.
- Kolenikov, Stanislav. 2014. Calibrating survey data using iterative proportional fitting (raking). *The Stata Journal*, 14(1): 22-59.
- Olson, Kristen. 2018. “Single imputation”. SOCI902: Analysis of Complex Survey Data course notes. University of Nebraska-Lincoln.
- Olson, Kristen. 2018. “Weights”. SOCI902: Analysis of Complex Survey Data course notes. University of Nebraska-Lincoln.
- Olson, Kristen. 2018. “Weights and software”. SOCI902: Analysis of Complex Survey Data course notes. University of Nebraska-Lincoln.