

2022-2023 National Survey of Family Growth: Weighting Methodology

Prepared Under Contract #75D30120C09732 with
National Center for Health Statistics
3311 Toledo Road
Hyattsville, MD 20782

Prepared by:
Taylor Lewis, Stephanie Zimmer, Jennifer Cooney, Andy Peytchev
RTI International
3040 East Cornwallis Road
Research Triangle Park, NC 27709-2194

August 2025

Table of Contents

Executive Summary.....	1
1. Introduction	1
2. Developing Base Weights.....	2
2.1 First Stage: Selection of PSUs.....	2
2.2 Second Stage: Selection of SSUs	3
2.3 Third Stage: Selection of Housing Units.....	3
2.4 Fourth Stage: Selection of Individuals	4
2.5 Fifth Stage: Double-Sampling for Nonresponse.....	5
3. Address-Level Weight Adjustments	6
4. Individual-Level Weight Adjustments	7
Appendix: Procedure for Estimating the Age-Eligible Population Living in College Group Quarters.....	10
References	11

Executive Summary

The National Survey of Family Growth (NSFG) is a survey on fertility, family formation and change, family planning, reproductive health, and closely related topics. This survey has been conducted by the National Center for Health Statistics (NCHS) since 1973 and is a principal source of national estimates on a variety of fertility and family topics. The target population for the NSFG consists of all non-institutionalized women and men aged 15-49 as of first contact for the survey, living in households, and whose usual place of residence is the 50 United States or the District of Columbia.

This document summarizes the steps taken to develop final analysis weights for the 9,957 completed main surveys (5,586 females and 4,371 males) obtained over the 2-year data collection between January 2022 and December 2023. Specifically, this report documents the development of base weights to account for variable selection probabilities of households and individuals, nonresponse adjustments at both the household and individual level, and the calibration of nonresponse-adjusted weights to population benchmarks derived from the U.S. Census Bureau's American Community Survey (ACS).

1. Introduction

The purpose of this document is to detail methods undertaken to create a set of analysis weights enabling population-level estimation for the 5,586 female and 4,371 male respondents to the 2022-2023 NSFG. NSFG respondents were selected randomly from the target population of non-institutionalized individuals aged 15-49 living in the United States by way of a multistage clustered sample design. In addition to households containing one or more individuals in that age range, the NSFG target population includes age-eligible military personnel living off-base and individuals living in college group quarters (e.g., dormitories, fraternity/sorority houses) who are eligible to be sampled through the households of their parents/guardians.

Not all individuals have the same probability of selection into the NSFG. Females, Black individuals, and teens are sampled at higher rates than others. The first step in creating analysis weights for use in constructing unbiased estimates of target population quantities is to account for these differential sampling rates by assigning a *base weight*, defined as the inverse of the product of the selection probabilities across all five stages of sampling. In Section 2 of this document, we introduce notation and formulas for how these base weights were defined.

If every sampled household completed the screening survey, and each age-eligible sampled individual complied with the request to participate in the main survey, the base weights would be all that are necessary to produce unbiased estimates of target population attributes. This does not happen in practice. *Unit nonresponse* refers to when the sampled individual (or household) is either unable to be contacted, or upon contact refuses or is physically unable to participate. A key feature of the NSFG sample design that attempts to mitigate the potentially harmful effects of unit nonresponse is the double-sampling approach (Deming, 1953; Hansen & Hurwitz, 1946) discussed in Section 7 of the **Sample Design** report. This was operationalized in NSFG by subsampling nonrespondents after a fixed portion of the quarterly data collection period and designating these cases for a targeted follow-up attempt, with a higher promised incentive, during the latter portion of the data collection period. Base weights of the subsampled nonrespondents were inflated to represent those not subsampled. Further adjustments for unit nonresponse are discussed in Section 3.

Lastly, Section 4 of this document discusses the third step in weighting, during which nonresponse adjusted weights are calibrated to population benchmarks derived from the U.S. Census Bureau. The ACS 1-year Public Use Microdata Sample (PUMS) file was used for this purpose. The ACS 1-year

microdata file is released annually, usually in December, to represent population figures for the year prior. Weights for the 2022-2023 NSFG data collection used the file released in late 2023 reflecting population figures as of calendar year 2022. Note that the ACS microdata include separate figures for people living in households, non-institutionalized group quarters, and institutionalized group quarters. Individuals living in college group quarters are part of the target population but represent only a portion of the ACS figures for those living in non-institutionalized group quarters. A procedure for estimating the share of the non-institutionalized group quarters population living in college group quarters, so that they could be included in the target population totals, is described in the Appendix.

2. Developing Base Weights

This section defines notation and formulas to express the probability of an individual being selected for the NSFG as a function of five subsequent stages of sampling. More details on the composition and descriptive characteristics of sampling units at the various sampling stages are provided in the **Sample Design** report.

2.1 First Stage: Selection of PSUs

The first stage of sampling was to select *primary sampling units* (PSUs), which consist of individual counties or groupings of contiguous counties. Note that one sample of 222 PSUs was drawn in advance of the 8-year data collection period spanning calendar year 2022 through 2029, and PSUs were subsequently allocated to one or more years; see Table 3 of the **Sample Design** report. Prior to this selection, PSUs were stratified into 8 groupings based on the cross-classification of the four Census Bureau regions and an indicator of the PSU being a Metropolitan Statistical Area. Additionally, there were three separate strata of 22 certainty PSUs, 9 of which are slated to be in the sample all 8 years (“Certainty Group 1”), and the other 13 of which appear in 4 adjacent years over the 8-year period (“Certainty Group 2a” and “Certainty Group 2b”).

The probability of selection for PSU i in design stratum h is

$$\pi_{hi}^0 = n_h \frac{M_{hi}}{M_h}$$

where n_h is the number of PSUs sampled from design stratum h and

$$M_{hi} = 2.6 * POP_HH_BLK_1459_{hi} + POP_HH_NONBLK_1459_{hi}$$

where $POP_HH_BLK_1459_{hi}$ is the estimated count of individuals in the target population in the PSU who are Black and $POP_HH_NONBLK_1459_{hi}$ is the same for those who are not Black. Details on how the population estimates are calculated are included in the Appendix of the **Sample Design** report.

Note that $M_h = \sum_{i=1}^{N_h} M_{hi}$ is the sum measure of size (MOS) across the population of N_h PSUs in design stratum h . A few exceptions were made if $\frac{M_{hi}}{M_h} > \frac{1}{n_h}$. In those instances, $\pi_{hi}^0 = 1$, which is also how selection probabilities have been assigned for certainty PSUs.

The last component of the first-stage selection probability is an allocation factor, a_{hi} , attributable to the fact that not all PSUs are released (i.e., fielded) in a single year. For PSUs in Certainty Group 1, $a_{hi} = 1$. For PSUs designated for Certainty Group 2a (Years 1 – 4), $a_{hi} = \frac{6}{13}$. For PSUs randomly designated in Certainty Group 2b (Years 5 – 8), $a_{hi} = \frac{7}{13}$. Lastly, for non-certainty PSUs, $a_{hi} = \frac{25}{200} = \frac{1}{8}$.

The final probability of selection for a PSU in a given year is then defined as

$$\pi_{hi} = a_{hi} \pi_{hi}^0$$

2.2 Second Stage: Selection of SSUs

The second stage of sampling was to select *secondary sampling units* (SSUs), also commonly referred to as *segments*, which are defined as Census Block Groups, or adjacent groupings thereof, within a PSU. Let M_{hij} represent the weighted MOS for SSU j in PSU i in design stratum h . Specifically,

$$M_{hij} = 2.6 * POP_HH_BLK_1459_{hij} + POP_HH_NONBLK_1459_{hij}$$

where $POP_HH_BLK_1459_{hij}$ is the estimated count of individuals in the target population in the SSU who are Black and $POP_HH_NONBLK_1459_{hij}$ is the same for those who are not Black.

The probability of selection for SSU j in PSU i in design stratum h is

$$\pi_{hij}^0 = n_{hi} \frac{M_{hij}}{M_{hi}}$$

where n_{hi} is the number of SSUs sampled in the PSU, and $M_{hi} = \sum_{j=1}^{N_{hi}} M_{hij}$ is the sum MOS for all N_{hi} SSUs in the PSU.

The last component of the second-stage selection probability is an allocation factor, a_{hij} , attributable to the fact that not all selected SSUs are fielded a single year. For SSUs in PSUs in Certainty Group 1, $a_{hij} = \frac{1}{8}$. For SSUs in PSUs in Certainty Groups 2a and 2b, $a_{hij} = \frac{1}{4}$. For SSUs in all other PSUs, $a_{hij} = 1$.

The final probability of selection for an SSU in a given year is then defined as

$$\pi_{hij} = a_{hij} \pi_{hij}^0$$

2.3 Third Stage: Selection of Housing Units

The third stage of sampling is to select individual addresses, which serve as proxies for housing units (or households). As discussed in Section 5 of the **Sample Design** report, there were two possible sources used to compile the sampling frame of addresses within an SSU: (1) RTI's address-based sampling (ABS) frame or (2) addresses listed from a field enumeration operation in SSUs when the estimated net coverage rate (Harter et al., 2021) of the ABS frame was less than 85%. Of the 979 SSUs sampled during the 2022-2023 data collection period, only 69 (about 7%) were listed. Lists of addresses for the remaining 93% of SSUs came directly from the ABS frame. In certain scenarios where an SSU designated to be listed was expected to have an extremely large number of addresses or an expansive land area to canvass, it was sub-segmented into a very small number of more manageable areas—most often two or three—and one was selected at random to be listed.

Addresses derived from the ABS frame were stratified based on the predicted likelihood that the address contained at least one age-eligible individual—see Table 4 of the **Sample Design** report. Three approximately equally sized strata were created based on the ranked likelihoods (low, medium, and high), and addresses from the higher likelihood strata were sampled at a higher rate. This was done to improve data collection efficiency, as any household without at least one individual aged 15-49 was considered ineligible and screened out of the survey. Note that addresses from listed SSUs constituted a fourth stratum.

The probability of selecting address l in SSU j in PSU i in design stratum h is a function of its SSU-level address stratum, k (i.e., one of three age-eligibility strata or the stratum of listed addresses), and a sub-segmentation factor c_{hijk} where, aside from a few special cases of SSUs deemed too large to list in their entirety, $c_{hijk} = 1$. If we let N_{hijk} denote the number of addresses in stratum k within the given SSU, and n_{hijk} the number of addresses sampled released, the address selection probability is

$$\pi_{hijkl} = \frac{1}{c_{hijk}} \frac{n_{hijk}}{N_{hijk}}$$

The initial base weight for a sampled address was defined as the inverse of the product of the first three stages' selection probabilities, or

$$w_{hijkl} = (\pi_{hi} \times \pi_{hij} \times \pi_{hijkl})^{-1}$$

An address's base weight can be interpreted as the number of addresses in the population that the sampled address represents. Address-level base weights sum to around 35.5 million addresses per quarter, or about 142 million addresses per year.

2.4 Fourth Stage: Selection of Individuals

The fourth stage of sampling was to select a single individual from the roster of all age-eligible individuals in the sampled address. The first step in this process was to request an adult household member to complete a short screener. The screener determined whether there were one or more individuals between the age of 15 and 49 living at the address. If no age-eligible individuals were present, the informant was thanked for their time, the address was coded as ineligible, and data collection was discontinued. For addresses with one or more age-eligible individuals present, the screener populated a roster of individual names (or initials) and a small number of critical demographics, such as age, sex, and race/Hispanic origin. From this roster, one individual was selected at random for the main survey using a pre-programmed feature built into the screener software.

The selection of an individual was not performed using a simple random sampling approach. Rather, one individual was selected with probability proportional to size according to the MOS figures in Table 1. In both 2022 and 2023, the MOS was contingent on an individual's age and sex. Females and teens were assigned a larger MOS to facilitate oversampling them relative to males and non-teens because NSFG targets 55% of all main survey completes to be female and 18.2% of all main survey completes to be teens. The MOS values used in 2022 were carried forward from those used in the 2017-2019 NSFG design, per Table 8 of <https://www.cdc.gov/nchs/data/nsfg/NSFG-2017-2019-Sample-Design-Documentation-508.pdf>. However, these MOS figures were decreased for males in 2023, thereby increasing the selection probabilities of females, to accommodate additional funding to increase the number of female completes by 500 in 2023.

Table 1: Within-household measures of size in 2022 and 2023

Category	2022	2023
Females 15-19	1.00	1.00
Females 20-49	0.25	0.25
Males 15-19	0.91	0.64
Males 20-49	0.23	0.16

For age-eligible individual m living at address l of stratum k within SSU j in PSU i in the primary design stratum h , we assigned an MOS according to Table 1. The selection probability of that individual was defined as

$$\pi_{hijklm} = \frac{F_{hijklm}}{F_{hijkl}}$$

where F_{hijkl} is the sum of the MOS values for all age-eligible individuals in the housing unit.

The initial base weight for sampled individual m , conditional on the household screening in as eligible, was defined as the inverse this selection probability, or

$$w_{hijklm} = \frac{1}{\pi_{hijklm}}$$

2.5 Fifth Stage: Double-Sampling for Nonresponse

The fifth and final stage of sampling is a nonrespondent follow-up technique referred to in the literature as *double sampling* or *two-phase sampling* for nonresponse (Deming, 1953; Hansen & Hurwitz, 1946), as used in the 2006-2010 and 2011-2019 NSFG. During the 2022-2023 NSFG, a random subsample of nonrespondents was selected to be offered an additional \$40 incentive to participate during the final 4 weeks of the 16-week multimode quarterly data collection administration. Because the quarterly data collection protocol during the 2022-2023 data collection period included three phases (as compared to two phases during the 2011-2019 period), this double sampling and enhanced protocol is referred to here as Phase 3. For more information on the quarterly data collection protocol, see the **Sample Design** report.

Prior to the onset of Phase 3, nonrespondents were stratified into one of six mutually exclusive statuses based on the cross-classification of breakoff status (yes or no) and data collection stage of nonresponse (screener, male, or female). Any breakoff was selected into the Phase 3 sample with certainty (i.e., assigned a selection probability of 1), whereas other strata were selected via simple random sampling at a rate between 40% and 50%, depending upon quarter. See discussion around Table 7 of the **Study Design and Data Collection Procedures** report for more details.

The probability of selecting a case is a function of its status at the onset of Phase 3. There was no change to the probability of selection for cases that responded prior to Phase 3. On the other hand, nonrespondents were assigned a modified probability of selection based on their status, which could be either at the address level (e.g., nonrespondent/breakoff at screener) or at the individual level (e.g., nonrespondent/breakoff at main survey). If we let $N_{hijkl'}$ denote the number of address-level cases eligible for Phase 3 and $n_{hijkl'}$ denote the number selected for a given quarter, address-level nonrespondents selected into Phase 3 had their probability of selection modified to

$$\pi_{hijkl'} = \pi_{hijkl} * \frac{n_{hijkl'}}{N_{hijkl'}}$$

Note that if a screener was completed during Phase 3, this probability was further multiplied by π_{hijklm} , calculated as described in the previous subsection of this report.

If we let $N_{hijklm'}$ denote the number of individual-level cases eligible for Phase 3 and $n_{hijklm'}$ denote the number selected for a given quarter, individual-level nonrespondents selected into Phase 3 had their probability of selection modified to

$$\pi_{hijklm'} = \pi_{hijklm} * \frac{n_{hijklm'}}{N_{hijklm'}}$$

Address- and individual-level base weights were modified to account for the double-sampling in Phase 3. Specifically, a screener nonrespondent's address-level base weight was multiplied by the inverse of its Phase 3 selection rate. Similarly, main survey nonrespondents selected for Phase 3 had their individual-level base weight multiplied by the inverse of their selection rate. Note that the base weight of any cases eligible for Phase 3, but not selected, was set to 0.

3. Address-Level Weight Adjustments

Starting with the address-level base weights, two adjustments were made to compensate for unit nonresponse at the household level. The first was an adjustment for unknown eligibility of nonresponding addresses. The two most common reasons for ineligibility at the address level were (1) vacancy and (2) a household not containing at least one individual between the ages of 15 and 49 at the time the screener was completed. However, eligibility status could only be determined for a portion of sampled addresses: those where the screener was completed and those where the data collection team received definitive evidence of vacancy (e.g., multiple undeliverable mailings sent back by USPS) or a field interviewer (FI) observing the building at a sampled address has been demolished or destroyed. Other (rare) examples of definitively known ineligibility was when an FI observed a sampled address resided within a military base or as part of a senior-living community. Aside from any of these scenarios, the remaining portion of addresses had unknown eligibility.

Unknown eligibility was handled by estimating the base-weighted eligibility rate within each address stratum in a PSU—that is, aggregating across all SSUs sampled in the PSU—and then multiplying the base weights of addresses with known eligibility by this rate. This adjustment assumes that the rate of eligibility in addresses where eligibility status could not be determined is equivalent to the rate of eligibility in the addresses where eligibility status could be determined, controlling on PSU and address stratum.

Using the adjusted address-level base weights produced by the first step, the second adjustment aimed to compensate for unit nonresponse by shifting weights of nonresponding addresses to similar responding addresses. This occurred within classes, or cells, populated with addresses having approximately equivalent (estimated) response propensities. In the terminology of Little and Rubin (2019), there is a *missing at random* (MAR) assumption for the sample as a whole but a *missing completely at random* (MCAR) assumption within each class.

Although a variety of techniques are used in practice to form classes, in the 2022-2023 NSFG, the Chi-Square Automatic Interaction Detector (CHAID) algorithm (Kass, 1980) built into PROC HPSPLIT in SAS® (SAS Institute Inc., 2015) was used. This algorithm fits an implicit model using a classification and regression tree approach (Breiman et al., 1984). It exploits a pool of potential covariates to recursively partition a data set into groupings referred to as *nodes*, or *leaves*, which are treated as weight adjustment classes for our purposes. These classifications are created by making a hierarchical sequence of binary splits that best explain residual variation in the outcome variable—in this case, a binary response indicator.

The algorithm commenced with several dozen covariates for consideration, including the adjusted base weight itself, PSU and SSU identifiers, address-level auxiliary variables appended to the address from commercial vendors, and summary measures at the Census Tract and Census Block Group levels derived from the ACS (e.g., percent of renter-occupied households, the median home value, the percent of individuals without health insurance, the percent of individuals living below the poverty level). The algorithm created 21 classes ranging from an unweighted response rate of 23.0% to 56.0%. Note that this rate treats a completed screener without any age-eligible individuals as a response. Using the 21 classes, respondent weights previously adjusted for unknown eligibility were multiplied by the inverse of the weighted response rate, per the recommendation of Kott (2012). The net result upon completing this step was that the weights of responding addresses were inflated to represent not only themselves, but the nonresponding addresses as well. The weights of nonresponding addresses were then set to 0, effectively dropping them from any subsequent weighting adjustments.

4. Individual-Level Weight Adjustments

Individual-level base weights calculated for households completing the screener were adjusted for main survey nonresponse in two subsequent steps. The first step was to adjust for main survey unit nonresponse by exploiting covariates known for both main survey respondents and nonrespondents, such as household information captured from the screener (e.g., demographics and household composition information, whether the screener informant was the individual selected for the main survey) or other information derived from RTI’s enhanced ABS frame. Analogously to the address-level unit nonresponse adjustment procedure, a battery of these potential covariates was supplied to the CHAID algorithm built into PROC HPSPLIT. A total of 36 classes were formed, and the weighted response rate within the class was used to inflate the weights of main survey respondents to compensate for individuals selected but who never responded.

The second weighting adjustment applied to the individual-level weights was a calibration step. This marked the final step in the weighting process. Using the nonresponse-adjusted individual-level weights from the prior step, the general exponential model approach (Folsom & Singh, 2000) built into SUDAAN’s WTADJUST (RTI International, 2012) procedure to altered individual-level weights such that totals broken out by categorizations of sex, age, race/Hispanic origin, highest education level attained, and marital status matched population figures estimated from the ACS 2022 PUMS File. In addition to maintaining univariate distributions, the calibration step ensured that the distributions of all two-way interactions of sex, age, race/Hispanic origin, and highest education level attained were maintained, as were the three-way interaction of sex, age, and race/Hispanic origin. Table 2 summarizes the source variables and corresponding categories that were controlled for as part of the calibration step.

Table 2. Variables and categories used for the calibration of individual-level weights.

Description	Source Variable(s)	Categories for Calibration
Sex	(Instrument)	1 = Male 2 = Female
Race/Hispanic origin	HISPRACE2	1 = non-Hispanic White 2 = non-Hispanic Black 3 = Hispanic 4 = Other
Age	AGER	1 = Less than 20 2 = 20 to 24 3 = 25 to 29 4 = 30 to 34 5 = 35 to 39 6 = 40 to 44 7 = 45 to 49
Marital Status ¹	MARSTAT/LMARSTAT	1 = Married 2 = Previously married 3 = Never married
Education	HIEDUC	1 = HS or less 2 = More than HS

¹ A respondent is considered married when MARSTAT (marital status) is coded 1. Else, if LMARSTAT (legal marital status) is coded between 3 and 5, then the respondent is considered previously married. All other respondents are considered never married.

Prior to finalizing the individual-level calibration step, a modest amount of weight trimming was conducted. Specifically, weights were trimmed such that the minimum weight would be no less than the 2nd percentile of the originally calibrated weight and the maximum weight be no more than the 98th percentile. As Table 3 shows, this trimming procedure did not substantively change key estimates, but it helped limit the precision loss attributable to the overall unequal weighting effect (UWE) (Kish, 1992), as seen from the reduced standard errors and design effects. Note that the final, calibrated, trimmed weight on the 2022-2023 NSFG public-use data files is named WGT2022_2023. For females, this weight ranges from 1,170.96 to 78,770.13, and sums to 74,936,918 individuals in the target population. For males, this weight ranges from 1,514.72 to 95,539.72, and sums to 75,700,206 individuals in the target population. Combined, the female and male weights sum to 150,637,124 individuals.

Table 3: Comparison of point estimates, standard errors, and design effects for the untrimmed calibrated weight and the trimmed weight

Outcomes	Female					
	Trimmed Weight			Untrimmed Weight		
	Percent	SE Percent	Design Effect	Percent	SE Percent	Design Effect
Age at first sex (<15)	14.5	0.77	2.20	14.4	0.78	2.27
Age at first sex (15-17)	39.6	1.20	2.73	39.9	1.35	3.34
Age at first sex (18+)	45.9	1.50	4.10	45.7	1.61	4.73
Ever cohabited	51.0	1.00	2.22	50.1	1.19	3.19
No live births	50.5	1.15	2.97	50.6	1.36	4.16
One live birth	15.3	0.58	1.51	14.7	0.65	1.89
Two or more live births	34.2	1.08	2.86	34.7	1.20	3.55
Intend a/another birth	46.2	0.83	1.51	45.9	0.98	2.09
Used contraception at first sex	68.9	1.05	2.30	69.4	1.15	2.75
Had sex in the last 12 months	83.7	0.64	1.36	83.6	0.72	1.71
Ever smoked at least 100 cigarettes	20.0	0.93	3.00	19.7	1.03	3.71
Ever had an HIV test outside of blood donation	43.1	1.18	3.16	42.4	1.34	4.03
Health care coverage in last 12 months	89.5	0.71	2.89	89.2	0.82	3.87
Received public assistance in the last 12 months	6.0	0.61	3.34	6.3	0.76	5.09
Ever pregnant	54.7	1.16	3.02	54.5	1.35	4.09

Outcomes	Male					
	Trimmed Weight			Untrimmed Weight		
	Percent	SE Percent	Design Effect	Percent	SE Percent	Design Effect
Age at first sex (<15)	18.1	1.10	2.12	18.3	1.29	2.93
Age at first sex (15-17)	42.7	1.08	1.26	42.8	1.18	1.49
Age at first sex (18+)	39.2	1.39	2.14	38.9	1.57	2.72
Ever cohabited	27.7	0.96	2.00	27.9	1.08	2.55
No biological children	50.0	1.22	1.97	48.7	1.32	2.31
One biological child	16.1	0.75	1.40	16.1	1.00	2.48
Two or more biological children	33.7	1.31	2.49	35.1	1.52	3.40
Intend a/another birth	55.1	0.88	1.32	55.5	0.96	1.58
Used contraception at first sex	74.1	1.29	2.28	74.8	1.51	3.12
Had sex in the last 12 months	84.7	0.75	1.44	84.6	1.01	2.64
Ever smoked at least 100 cigarettes	28.5	1.13	2.74	28.1	1.42	4.30
Ever had an HIV test outside of blood donation	33.3	1.08	2.29	32.9	1.24	2.99
Health care coverage in last 12 months	85.2	0.93	2.89	85.1	1.12	4.20
Received public assistance in the last 12 months	5.1	0.45	1.71	5.1	0.57	2.70

Appendix: Procedure for Estimating the Age-Eligible Population Living in College Group Quarters

The ACS microdata includes people living in households, non-institutionalized group quarters, and institutionalized group quarters. Individuals living in college group quarters are only a portion of those living in non-institutionalized group quarters, but no separate indicator is provided in the microdata to differentiate them. This appendix describes a three-step procedure to estimate this population by combining information from ACS microdata with supplementary ACS tabular data. The three steps are as follows:

1. Estimate number of people living in non-institutionalized group quarters by sex, age, and race/Hispanic origin using ACS microdata.
2. Estimate the proportion of people living in college dormitories among those in non-institutionalized group quarters by sex, age, and race/Hispanic origin using ACS tabular data.
3. Multiply the estimated proportions in Step 2 by the totals in Step 1 to get an estimate of number of people living in college dorms by age, sex, and race/Hispanic origin.

More details on each step are provided below.

Step 1: Using the ACS microdata, the number of people living in households and non-institutionalized group quarters is tabulated by sex, age, and race/Hispanic origin using the analysis weights provided in the file. The age groupings used for this calculation are 15-17, 18-19, 20-24, 25-29, 30-34, 35-39, 40-44, and 45-49 and the race/Hispanic origin groupings used are Hispanic, non-Hispanic Black, and non-Hispanic other.

Step 2: The Census Bureau also publishes ACS data via tables that contain the marginal numbers of people living in college dormitories by sex, age, and race/Hispanic origin from the ACS with age groupings of 15-17, 18-24, 25-34, 35-44, and 45-54. Using data from these tables, raking (Kalton & Flores-Cervantes, 2003) can be used to estimate the number of people living in college group quarters for the joint distribution of these covariates. Using these estimated totals, the proportion of people living in college group quarters among those in non-institutionalized group quarters can be estimated.

Step 3: The estimated proportions of individuals living in college group quarters estimated in Step 2 are multiplied by the counts of individuals living in non-institutionalized group quarters to get an estimate of this target population domain. The age groups can then be collapsed down to the poststratification cells of individuals aged 15-19, 20-24, 25-29, 30-34, 35-39, 40-44, and 45-49.

References

- Breiman, L., Friedman, J., Olshen, R., & Stone, C. (1984). *Classification and regression trees*. Belmont, CA: Wadsworth, Inc.
- Deming, W. (1953). On a probability mechanism to attain an economic balance between the resultant error of nonresponse and the bias of nonresponse. *Journal of the American Statistical Association*, *48*, 743-772.
- Folsom, R., & Singh, A. (2000). The generalized exponential model for sampling weight calibration for extreme values, nonresponse, and poststratification. *Proceedings of the Joint Statistical Meetings*. Alexandria, VA: American Statistical Association.
- Hansen, M., & Hurwitz, W. (1946). The problem of nonresponse in sample surveys. *Journal of the American Statistical Association*, *41*, 517-529.
- Harter, R., Morton, K., Amaya, A., & Brown, D. (2021). Estimating net coverage of segments in ABS frames. *Field Methods*, *33*, 68-84.
- Kalton, G., & Flores-Cervantes, I. (2003). Weighting methods. *Journal of Official Statistics*, *19*, 81-97.
- Kass, G. V. (1980). An exploratory technique for investigating large quantities of categorical data. *Applied Statistics*, *29*, 119-127.
- Kish, L. (1992). Weighting for unequal Pi. *Journal of Official Statistics*, *8*, 183-200.
- Kott, P. (2012). Why one should incorporate the design weights when adjusting for unit nonresponse using response homogeneity groups. *Survey Methodology*, *38*, 95-99.
- Little, R., & Rubin, D. (2019). *Statistical analysis with missing data*. 3rd ed. New York: Wiley.
- RTI International. (2012). *SUDAAN: Statistical software for weighting, imputing, and analyzing data, Release 11*. Research Triangle Park, NC: Research Triangle Institute.
- SAS Institute Inc. (2015). *SAS/STAT® User's guide*. Cary, NC: SAS Institute Inc.