

COMPARABILITY OF DATA: BRFSS 2001

The BRFSS is a cross-sectional surveillance survey currently involving 54 reporting areas. It is important to note that any survey will have natural variation over sample sites; therefore, some variation between states is to be expected. The complex sample design and the multiple reporting areas complicate the analysis of the BRFSS. Although CDC works with the states to minimize deviations, in 2001 there were some deviations in sampling and weighting protocols, sample size, response rates, and collection or processing procedures. In addition, California's questionnaire had a few minor differences in wording of question. The following section identifies other known variations for the 2001 data year.

A. 2001 Data Anomalies and Deviations from Sampling Frame and Weighting Protocols

In several states, a portion of sample records intended for use during one month were completed in another month. This deviation will disproportionately affect analyses based on monthly, rather than annual data.

Several states did not collect data for all 12 months of the year or completed interviews in calendar year 2002. Florida, Hawaii, and Texas did not report any interviews in January. New Jersey and the Virgin Islands did not complete any interviews in January or February. The District of Columbia did not complete any interviews in January, February, or March. Nevada completed 302 interviews in January or February, 2002. Illinois, New Mexico, South Dakota, and Wisconsin had small numbers of completed interviews in January, 2002.

More information about the quality of the survey data can be found in the *2001 BRFSS Summary Data Quality Report*.

B. Other 2001 limitations of the data

Telephone coverage varies by state and also by sub-population. Telephone coverage averages 94.5% for U.S. states as a whole, but ranges from 1.8% noncoverage in Delaware, to 13.3% in New Mexico. It is estimated that 10% of households in Puerto Rico are without telephones.

Illinois used a dual questionnaire and collected data on core items involving immunization, cholesterol awareness, hypertension awareness, arthritis, firearms, disability, physical activity, prostate cancer screening, colorectal cancer screening, and HIV/AIDS knowledge and prevention, and a module on fruit and vegetable consumption on approximately half of respondents. Fruit and vegetable consumption questions were not asked in August but were asked of all respondents in December. Questions on hypertension awareness, cholesterol awareness, arthritis,

immunization, firearms, disability, physical activity, prostate cancer screening, colorectal cancer screening, and HIV/AIDS knowledge and prevention were asked of all eligible respondents in August.

California modified the wording of the health plan, diabetes, and the frequency of alcohol consumption questions. These questions may have limited comparability to those of other reporting areas.

Several states that used the tobacco indicators module inappropriately skipped question 6 (HOUSESMK) in the module for part of the year. The HOUSESMK field was left blank in these states if the question was not asked. The Tobacco Indicators module has a second problem where several states skipped question 4 (GETCARE) when question 3 (LASTSMK) had a response of 1, 2, 3, or 4. In some states the affected records were coded, 9 Refused, even though the respondents were not asked the questions, other states left the field blank.

The number of affected records by state and months for question 6, HOUSESMK are

State	Month(s)	Frequency
AK	January–February	177
KY	January	249
LA	January–February	370
MS	January–February	192
MT	January	123
NE	January–September	1263
OK	January–March	272
SD	February–March	243
WV	January–February	206

The number of affected records by state and months for question 4, GETCARE are

State	Month(s)	Frequency
AK	January	6
AZ	January–December	74
AR	January–April	21
IN	January–March	22
KY	January	12
LA	January–February	22
ME	January–December	63
MS	January	7
MT	January–February	13
NE	February, June–September	30
NJ	January–December	131
OK	January–May	48
PA	January–April	19
RI	January–December	109
SC	January–October	56
SD	January–March	14
VA	January–December	65
WY	January–May	19

More information about survey item nonresponse can be found in the *2001 BRFSS Summary Data Quality Report*.

STATISTICAL AND ANALYTIC ISSUES

Estimation Procedures

Unweighted data on the BRFSS are the actual responses of each respondent. Unweighted data represent results before any adjustment is made for variation in respondents= probability of selection, for disproportionate selection of population subgroups relative to the state=s population distribution, or nonresponse. Weighted BRFSS data represent results that have been adjusted to compensate for such differences. Irrespective of state sample design, use of the final weight in analysis is necessary if generalizations are to be made from the sample to the population.

Statistical Issues

The procedures for estimating variances given in most statistical texts and the programs available in most statistical software packages are based on the assumption of simple random sampling (SRS). The data collected in the BRFSS are obtained through a complex sample design; therefore, the direct application of standard statistical analysis methods for variance estimation and hypothesis testing may yield misleading results. There are computer programs available that take such complex sample designs into account. SAS Version 8's SURVEYMEANS and SURVEYREG procedures, SUDAAN, and EpiInfo=s C-Sample are among those suitable for analyzing BRFSS data. SAS and SUDAAN can be used for both tabular and regression analyses; SUDAAN has more available options. EpiInfo=s C-sample can be used to calculate simple frequencies and two-way cross-tabulations. These software products require knowing the stratum, the primary sampling units, and the record weight—all of which are on the master data file. For more information on calculating variance estimations using SAS, see SAS Institute, 1999 (10). For information about SUDAAN, see Shah, Barnwell, Bieler, 1997 (1). For information about EpiInfo, see Dean, et al, 1995 (2).

Although the overall number of persons in the BRFSS is quite large for statistical inference purposes, subgroup analyses can lead to estimators that are unreliable. Consequently, analysis of subgroups, especially within a single data year or geographic area, requires that the user pay particular attention to the subgroup sample size. Small sample sizes may produce unstable estimates. Reliability of an estimate depends on the actual **unweighted** number of respondents in a category, not on the weighted number. Interpreting and reporting weighted numbers that are based on a small, unweighted number of respondents can mislead the reader into believing that a given finding is much more precise than it actually is. The BRFSS follows a rule of not reporting or interpreting percentages based upon a denominator of fewer than 50 respondents (unweighted sample).

Analytic Issues

Advantages and Disadvantages of Telephone Surveys

Compared with in-person interviewing techniques, telephone interviews are easy to conduct and monitor, and cost efficient. However, telephone interviews have limitations. Telephone surveys may have higher levels of noncoverage than in-person interviews because a percentage of United States households cannot be reached by telephone. As mentioned earlier, approximately 98 percent of households in the United States have telephones. A number of studies have shown that the telephone and non-telephone populations are different with respect to demographic, economic, and health characteristics (3,4,5). Although the estimates of characteristics for the total population are unlikely to be substantially affected by the omission of the non-telephone households, some of the subpopulation estimates could be biased due to the noncoverage of households without telephones. Telephone coverage is lower for population subgroups such as blacks in the South, persons with low incomes, persons in rural areas, persons with less than 12 years education, persons in poor health, and heads of households under 25 years of age (6). However, post-stratification adjustments for age, race, and sex, and other weighting adjustments used for the BRFSS data minimize the impact of differences in noncoverage, undercoverage, and nonresponse at the State level. State-specific information on telephone coverage is available in the technical documentation section on www.cdc.gov/brfss.

Despite the above limitations, prevalence estimates from the BRFSS correspond well with findings from surveys based on in-person interviews, including studies conducted by the National Institute on Alcohol Abuse and Alcoholism, CDC's National Center for Health Statistics, and the American Heart Association (7). A summary of methodologic studies of BRFSS is provided in the publication section on www.cdc.gov/brfss.

Surveys based on self-reported information may be less accurate than those based on physical measurements. For example, respondents are known to under report weight. Although this type of potential bias is an element of both telephone and in-person interviews, the under reporting should be considered by the analyst interpreting self-reported data. When measuring change over time, this type of bias is likely to be constant, and therefore not a factor in trend analysis.

Aggregating Data Over Time

When data from one time period are insufficient for estimating the prevalence of a risk factor, data may be combined for several periods as long as the periods being combined are not times during which the prevalence of the risk factor

of interest has been substantially changing. One method that can be used to assess the stability of the prevalence estimates is discussed below (7).

1. Compute the prevalence for the risk factor for each period.
2. Rank the estimates from low to high.
3. Identify a statistical test appropriate for comparing the lowest and the highest estimates at the 5% level of significance. For example, depending on the type of data, a t-test or the sign test might be appropriate.
4. Test the hypothesis that prevalence is not changing by using a two-sided test in which the null hypothesis is that the prevalences are equal.
5. Determine whether the resulting difference could be expected to occur by chance alone less than 5% of the time (i.e., test at the 95% confidence level).

Analyzing Subgroups

When the prevalence of risk factors does not change rapidly over time, data combined for two or more years may provide a sufficient number of respondents so that additional prevalence estimates can be made for population groups (such as age/sex/race subgroups or county populations). Before combining data for subgroups, determine whether the total number of respondents will yield the precision needed. The level of precision needed depends upon the intended use of the estimate. For example, greater precision would be required to justify implementing expensive programs than that for general information only.

The table below shows the sample size required for each of several levels of precision based on a calculation in which the estimated risk factor prevalence is 50% and the design effect is 1.5.

<u>Precision Desired</u>	<u>Sample Size Needed</u>
2%	3600
4%	900
6%	400
8%	225
10%	144
15%	64
20%	36

Precision is indicated by the width of the 95% confidence interval around the prevalence estimate. For example, a desired precision of 2% means that the 95% confidence interval is + or - 2% of 50%, or 48–52%. As shown in the table, to yield this high a level of precision, the sample size required is about 3,600 persons. When a lower level of precision is acceptable, the sample size can be considerably smaller.

The design effect is a measure of the complexity of the sampling design and indicates how the design differs from simple random sampling. It is defined as the variance for the actual sampling design divided by the variance for a simple random sample of the same size (7, 8). For most risk factors in most states, the design effect is less than 1.5. If it is more than 1.5, however, sample sizes may need to be larger than those shown here.

The standard error of a percentage is largest at 50% and decreases as a percentage approaches 0% or 100%. From this perspective, the required sample sizes above are conservative estimates. They should be reasonably valid for percentages between 20% and 80% but may significantly overstate the required sample sizes for smaller or larger percentages.

As a cautionary note, users should remember that the reliability of an estimate depends on the actual, unweighted number of respondents in a category, not on the weighted number. Interpreting and reporting weighted numbers that are based on a small, unweighted number of respondents can mislead the reader into believing that a given finding is much more precise than it actually is. **The CDC strongly urges all users to follow the general rule of not reporting or interpreting percentages based upon a denominator with fewer than 50 unweighted respondents.**

Creating Synthetic Estimates

Sample sizes may still be inadequate for risk factor estimates for some geographic areas (i.e., counties) or subpopulations (i.e., persons with diabetes) even after combining data for several years. In such situations, the analyst may wish to derive synthetic estimates by extrapolating from the BRFSS data collected at the state level.

Synthetic estimates can be calculated by using the population estimates for the subgroup of interest and the state BRFSS risk factor prevalences for that subgroup. This approach assumes that the risk factor prevalences for specific subgroups in each area are the same as the statewide risk factor prevalences for the same subgroups. For example, it assumes that the risk factor prevalences for black women in every county of a state are the same as those for black women in the entire state. The accuracy of the estimate depends on the validity of this assumption, which is often impossible to judge. However, a ballpark estimate may be sufficient for establishing broad goals and objectives for prevention strategies. For a discussion of the precision of such estimates, see Levy and Lemeshow (9).

An example for estimating the number of persons with hypertension in a hypothetical county, as well as the overall prevalence of hypertension in that county is shown below. The sex and race distribution of the county's population differs from the statewide population, and these differences need to be taken into account. By developing a table like the one below, a synthetic estimate for the overall county prevalence of hypertension can be made.

Synthetic Estimates of Prevalence of Hypertension in a Hypothetical County, 1990

State Subgroup	Prevalence* 1990	County Population 1990	County Population with Hypertension 1990
<i>Men</i>			
White	15.6	10,000	1,560
Black	27.0	25,000	6,750
<i>Women</i>			
White	19.5	12,000	2,340
Black	26.5	28,000	7,420
<i>Total</i>		75,000	18,070

*Per 100 persons

The state prevalence values, given as rates per 100 persons, are computed from the BRFSS data. The estimated number of persons with hypertension for each race-sex group in the county was obtained by multiplying the statewide prevalence for that group by the county population for the group. To determine the total county prevalence, the number of people with hypertension in each race-sex group in the county were summed and this sum (18,070) was divided by the county's total population (75,000) to yield an overall prevalence of 24.1 per 100 persons.

Creating Direct Estimates

If the subpopulation sample size is sufficient to do so, analysts may choose to produce direct estimates. SUDAAN or a similar program will be needed for direct estimates. The subarea (i.e., county) is treated as a population domain for which the risk estimate is sought, and will be defined as a SUBGROUP variable in SUDAAN. Temporal and spatial stratification must be incorporated into the estimates of variable, by inclusion in the NEST statement in SUDAAN. If possible, it is desirable to re-adjust the poststratification weight (_POSTSTR) to the age-by-race-by-gender population distribution of the small area (i.e., county).

To locally post-stratify the CDC BRFSS weights used for the direct estimate, post-stratify _WT1 to the population of interest. The equivalent local final weight is a product of _WT1 and the local poststratification factor.

New Race Variables

Starting in 2001, the BRFSS allowed respondents to choose more than one race. This change required a revision of calculated race variables. This section describes the coding of the race questions asked on the BRFSS and the variables that were calculated from them. The variable names are those assigned on the BRFSS SAS data file and by the SASOUT.SAS program that creates a SAS data file from a BRFSS ASCII data file. The column numbers are the location of the variable on the BRFSS ASCII file.

The following race questions were asked on the 2001 BRFSS:

13.3. Which one or more of the following would you say is your race? (113–118)

		Please Read
Mark all that apply	1	White
	2	Black or African American
	3	Asian
	4	Native Hawaiian or Other Pacific Islander
	5	American Indian, Alaska Native
		or
	6	Other [specify]
	8	No additional choices
Do not read these responses	7	Don't know/Not sure
	9	Refused

If more than one response to Q13.3, continue. Otherwise, go to Q13.5

13.4. Which one of these groups would you say best represents your race? (119)

- 1 White
- 2 Black or African American
- 3 Asian
- 4 Native Hawaiian or Other Pacific Islander
- 5 American Indian, Alaska Native
- 6 Other [specify]
- 7 Don't know/Not sure
- 8 Multiracial But Preferred Race Not Asked
- 9 Refused

Six columns were allocated for Question 13.3 (MRACE). Each race mentioned was coded in sequential columns starting with column 113. In some states, the responses for all records are in numerically ascending order; in other states, some responses are not in ascending order. The sequence was terminated by a 7, 8, or 9 unless all 6 races, including Other, were selected.

Question 13.4 (ORACE2, column 119) did not contain a response category of "8 Multiracial, But Preferred Race Not Asked" in the questionnaire. This was added subsequently to account for those few cases where a respondent answered 13.3 with more than one race but 13.4 was not asked.

MRACEORG (columns 616–621) is MRACE with trailing 7s, 8s, and 9s in columns 617–621 stripped off.

MRACEASC (columns 622–627) is MRACEORG with responses in ascending order.

_PRACE (columns 628–629), or Preferred Race, is calculated from MRACE and ORACE2. The values of _PRACE are

- 01= White
- 02= Black or African American
- 03= Asian
- 04= Native Hawaiian or Other Pacific Islander
- 05= American Indian or Alaskan Native
- 06= Other Race
- 07= No Preferred Race
- 08= Multiracial But Preferred Race Not Asked
- 77= Don't know/Not sure
- 99= Refused

_PRACE equals MRACE if the respondent answered MRACE with only one race. If the respondent indicated more than one race, then _PRACE equals 8 if ORACE2 was not asked. If ORACE2 was asked and the respondent

indicated a preferred race (responses 1–6), then _PRACE equals ORACE2. If ORACE2 was asked and the respondent did not indicate a preferred race (responses 7 or 9), then _PRACE equals 07. If the respondent gave an answer of Don't know/Not sure (7) or Refused (9) to MRACE, _PRACE equals 77 or 99 respectively.

_MRACE (columns 786–787), or Multiple Race, is calculated from MRACE. The values of _MRACE are

- 01= White Only
- 02= Black or African American Only
- 03= Asian Only
- 04= Native Hawaiian or Other Pacific Islander Only
- 05= American Indian or Alaskan Native Only
- 06= Other Race Only
- 07= Multiracial
- 77= Don't know/Not sure
- 99= Refused

_MRACE equals MRACE if the respondent answered MRACE with only one race. If the respondent indicated more than one race, then _MRACE equals 7. If the respondent gave an answer of Don't know/Not sure (7) or Refused (9) to MRACE, _MRACE equals 77 or 99 respectively.

_CNRACE (column 723), or the Number of Census Race Categories Chosen, is calculated from MRACE. It is the number of columns with a value of 1–5. Its value can vary from 0–5.

_CNRACEC (column 724), or the Number of Census Race Categories Chosen, Collapsed, is collapsed from _CNRACE. _CNRACEC equals missing or blank when _CNRACE equals 0, 1 when _CNRACE equals 1, and 2 when _CNRACE equals 2–5.

The following variables also involve responses to the Hispanic or Latino origin question, HISPANC2. This question was

- 13.2. Are you Hispanic or Latino? (112)
- | | |
|---|---------------------|
| 1 | Yes |
| 2 | No |
| 7 | Don't know/Not sure |
| 9 | Refused |

RACE2 (column 720), which replaces RACE from previous years, is calculated from HISPANC2 and _MRACE. The values of RACE2 are

- 1= White Only, Non-Hispanic
- 2= Black Only, Non-Hispanic

3= Asian Only, Non-Hispanic

4= Native Hawaiian or Other Pacific Islander Only, Non-Hispanic

5= American Indian or Alaskan Native Only, Non-Hispanic

6= Other Race Only, Non-Hispanic

7= Multiracial, Non-Hispanic

8= Hispanic

9= (Don't know/Not sure or Refused Hispanic Origin) or (Not Hispanic and Don't know/Not sure or Refused Race) [HISPANC2 IN (7,9) OR (HISPANC2 EQ 2 AND _MRACE IN (77,99))]

Hispanics are assigned a code of 8 and the 77's and 99's from _MRACE are assigned a code of 9. Otherwise, RACE2 equals _MRACE.

_RACEGR2 (column 722), which replaces _RACEGR from previous years, is collapsed from RACE2. The values of _RACEGR2 are

1= White Only, Non-Hispanic

2= Black Only, Non-Hispanic

3= Other Race Only, Non-Hispanic

4= Multiracial, Non-Hispanic

5= Hispanic

9= (Don't know/Not sure or Refused Hispanic Origin) or (Not Hispanic and Don't know/Not sure or Refused Race)

_RACEG2 (column 721), which replaces _RACEG from previous years, is also collapsed from RACE2. The values of _RACEG2 are

1= White Only, Non-Hispanic

2= Non-White Only, Multiracial, or Hispanic

9= (Don't know/Not sure or Refused Hispanic Origin) or (Not Hispanic and Don't know/Not sure or Refused Race)

REFERENCES

1. Shah BV, Barnwell BG, Bieler GS. SUDAAN User's Manual, Release 7.5, Research Triangle Park, NC: Research Triangle Institute, 1997.
2. Dean AG, Dean JA, Coulombier D, Brendel KA, Smith DC, Burton AH, Dicker RC, Sullivan K, Fagan RF, Arner TG. Epi Info, Version 6.0: A word processing, database, and statistics program for public health on IBM-compatible microcomputers. Centers for Disease Control and Prevention. 1995.
3. Groves RM, Kahn RL. Surveys by Telephone: A national comparison with personal interviews, New York, Academic Press, 1979.
4. Banks MJ. Comparing health and medical care estimates of the phone and nonphone populations. Proceedings of the Section on Survey Research Methods, American Statistical Association, 1983, pp 569-574.
5. Thornberry OT, Massey JT. Trends in United States Telephone Coverage Across Time and Subgroups. In Groves, RM et al editors Telephone Survey Methodology, pp. 25-49, New York, John Wiley & Sons, 1988.
6. Massey JT, Botman SL. Weighting Adjustments for Random Digit Dialed Surveys. In Groves, RM et al editors Telephone Survey Methodology, pp. 143-160, New York, John Wiley & Sons, 1988.
7. Frazier EL, Franks AL, Sanderson LM. Behavioral Risk Factor Data. In Using chronic disease data: A handbook for public health practitioners, pp 4.1-1.17. Centers for Disease Control and Prevention. 1992.
8. Groves RM. Survey Errors and Survey Costs. New York: John Wiley and Sons, 1989; 265, 271-272.
9. Levy PS, Lemeshow S. Sampling of Populations: Methods and Applications. New York: John Wiley and Sons, 1991; 347-350.
10. SAS Institute Inc., SAS/STAT User's Guide, Version 8. Cary, NC: SAS Institute, Inc., 1999; 3181-3272.

