

**Notes On Weights, Produced by Knowledge Networks, Amended by the Stanford Research Team,  
Applicable to Version 2.0 and later versions of the data.**

## **Sample Weighting**

The design for a KnowledgePanel<sup>SM</sup> sample begins as an equal probability sample that is self-weighting with several enhancements incorporated to improve efficiency. Since any alteration in the selection process is a deviation from a pure equal probability sample design, statistical weighting adjustments are made to the data to offset known selection deviations. These adjustments are incorporated in the sample's **base weight**.

There are also several sources of survey error that are an inherent part of any survey process, such as non-coverage and non-response due to panel recruitment methods and to inevitable panel attrition. We address these sources of sampling and non-sampling error using a **panel demographic post-stratification weight** as an additional adjustment.

Lastly, a set of **study-specific post-stratification weights** are constructed for the How Couples Meet and Stay Together Project data to adjust for sample design and survey non-response.

## **The Base Weight**

In a Knowledge Networks panel sample, there are six known sources of deviation from an equal probability of selection design. These are corrected in the Base Weight and are described below.

### **1. Under-sampling of telephone numbers unmatched to a valid mailing address**

An address match is attempted on all the Random Digit Dial (RDD) generated telephone numbers in the sample after the sample has been purged of business and institutional numbers and screened for non-working numbers. The success rate for address matching is in the 60-70% range. The telephone numbers with valid addresses are sent an advance letter, notifying the household that they will be contacted by phone to join KnowledgePanel. The remaining, unmatched numbers are under-sampled as a recruitment efficiency strategy. Advance letters improve recruitment success rates. Under-sampling stopped between July 2005 and April 2007. It was resumed in May 2007 with a sampling rate of 0.75.

### **2. RDD selection proportional to the number of telephone landlines reaching the household**

As part of the field data collection operation, information is collected on the number of separate telephone landlines in each selected household. A multiple line household's selection probability is down weighted by the inverse of its number of landlines.

### **3. Some minor oversampling of Chicago and Los Angeles due to early pilot surveys**

Two pilot surveys carried out in Chicago and Los Angeles when the panel was first being built increased the relative size of the sample from these two cities. With natural attrition and growth in size, the impact is disappearing over time. It remains part of our base adjustment weighting because of a small number of extant panel members from that nascent panel cohort.

#### **4. Early oversampling the four largest states and central region states**

At the time when the panel was first being built, survey demand in the four largest states (California, New York, Florida, and Texas) required over-sampling during January-October 2000. Similarly, the central region states were over-sampled for a brief period. These now diminishing effects still remain in the panel membership and thus require weighting adjustments for these geographic areas.

#### **5. Under-sampling of households not covered by the MSN<sup>®</sup> TV service network**

Certain small areas of the U.S. are not serviced by MSN<sup>®</sup>, thus MSN<sup>®</sup>TV units cannot be used. We under-sample households in these areas and use other Internet Service Providers for their Internet access.

#### **6. Oversampling of African- American and Hispanic telephone exchanges**

As of October 2001, we began over-sampling telephone exchanges with a higher density of minority households (uniquely African American and Hispanic) to increase panel membership for those groups. These exchanges are oversampled at approximately twice the rate of other exchanges. This over-sampling is corrected in the base weight.

### **The Panel Demographic Post-stratification Weight**

Generally, to reduce the effects of any non-response and non-coverage bias in a recruited panel, a post-stratification adjustment is applied using demographic distributions from the most recent data from the Current Population Survey (CPS). The post-stratification variables would include age, race, gender, Hispanic ethnicity, and education plus an Internet adjustment based on KnowledgePanel recruitment data. This weighting adjustment would be applied prior to the selection of any client sample from KnowledgePanel and usually constitutes the starting weights for survey samples when they are entirely composed of active panel members. For the How Couples Meet and Stay Together Project, however, the sample of active panel members was augmented with a sample of withdrawn panel members in an attempt to maximize sample size for GLB respondents. Because the usual post-stratification step does not apply to withdrawn panel members (it is based on the “active” panel only), this step was not used in this study to allow for the blending of active and withdrawn members. Instead, the base weight was used as the starting weight, letting the final post-stratification procedure (see next section) make the necessary demographic adjustments. This is a typical and successful solution for sample designs such as the one used in this study.

### **The Final Post-Stratification Weights for the How Couples Meet and Stay Together Project**

Once all the How Couples Meet and Stay Together Project data were returned from the field, we proceeded with a post-stratification process to adjust for any survey non-response and also any non-coverage due to the study-specific sample design. Demographic and geographic distributions for the population ages 18+ who are GLB or non-GLB from KnowledgePanel were used as benchmarks in this adjustment.

The following benchmark distributions were utilized for this post-stratification adjustment:

- Gender (Male, Female) x GLB (Yes, No)
- Age (18-29, 30-44, 45-59, 60+) x GLB (Yes, No)
- Race/Hispanic ethnicity (White/Non-Hispanic, Black/Non-Hispanic, Other/Non-Hispanic, Hispanic, 2+ Races/Non-Hispanic) x GLB (Yes, No)
- Education (Less than High School, High School, Some College, Bachelor and higher) x GLB (Yes, No)

- Census Region (Northeast, Midwest, South, West) x GLB (Yes, No)
- Metropolitan Area (Yes, No) x GLB (Yes, No)
- Internet Access (Yes, No) x GLB (Yes, No)

Comparable distributions were calculated using all completed cases from the field data. The completed cases include 1) respondents from the general population sample (including both GLB and non-GLB panelists), 2) respondents from the GLB augmentation sample (including both current and withdrawn panelists) and 3) respondents who had declined to answer the GLB identification question on the internal profile survey, but who reported being GLB upon re-contact for this survey. The last two categories in combination make up the GLB augmentation sample. Since study sample sizes are typically too small to accommodate a complete cross-tabulation of all the survey variables with the benchmark variables, an iterative proportional fitting is used for the post-stratification weighting adjustment. This procedure adjusts the sample data back to the selected benchmark proportions. Through an iterative convergence process, the weighted sample data are optimally fitted to the marginal distributions.

After this final post-stratification adjustment, the distribution of the calculated weights were examined to identify and, if necessary, trim outliers at the extreme upper and lower tails of the weight distribution separately by GLB and non-GLB status. The post-stratified and trimmed weights of the total respondents were scaled so that the weighted data sum to the actual sample size of total screened respondents (WEIGHT1 with n=4,002), total qualified (i.e., coupled) respondents (WEIGHT2 with n=3,009), and total qualified respondents by GLB and non-GLB status (WEIGHT3 with n=692 for GLB and n=2,317 for non-GLB).

We calculated additional post-stratification weights for the general population and GLB augmentation samples separately.

The following benchmark distributions were utilized for the post-stratification adjustment:

- Gender (Male, Female)
- Age (18-29, 30-44, 45-59, 60+)
- Race/Hispanic ethnicity (White/Non-Hispanic, Black/Non-Hispanic, Other/Non-Hispanic, Hispanic, 2+ Races/Non-Hispanic)
- Education (Less than High School, High School, Some College, Bachelor and higher)
- Census Region (Northeast, Midwest, South, West)
- Metropolitan Area (Yes, No)
- Internet Access (Yes, No)
- GLB (Yes, No) [used in the general population sample only]

After trimming outliers at the extreme upper and lower tails of the weight distribution, the post-stratified and trimmed weights of the general population sample respondents were scaled so that the weighted data sum to the actual sample size of the general population sample (WEIGHT4 with n=3,138) and qualified respondents (WEIGHT5 with n=2,377). We repeated the same for the GLB augmentation sample respondents and the post-stratified and trimmed weights were scaled so that the weighted data sum to the actual sample size of the GLB augmentation sample (WEIGHT4 with n=864) and the sample of qualified respondents (WEIGHT5 with n=632).

Two additional weights for GLB respondents were calculated excluding GLB cases from the withdrawn sample and from the group who had earlier declined to answer the GLB identification question. These weights are for total screened on-panel GLB respondents (WEIGHT6 with n=816) and total qualified (i.e., coupled) on-panel GLB respondents (WEIGHT7 with n=588).

Finally, all the original weights were scaled to reflect the actual size of the U.S. population, based on CPS benchmarks.

Base sampling weights are also included in the file for reference.

### **Application of the Weights**

Which of the calculated weights should be used depends on the unit of interest for analysis. An application summary by weight appears below.

The Weights That are Most Useful and Most Broadly Applicable are **weight1** and **weight2**.

Weight 1: This weight should be used for analysis of all screened respondents across both the general population and GLB augmentations samples, both those in couples and those not in couples.

Weight 2: This weight should be used for analysis of all coupled respondents across both the general population and GLB augmentations samples.

The Following weights are for more specialized circumstances, and therefore are labeled as “supplementary” weights in the public dataset:

Weight 3: This weight should be used for analysis of coupled respondents separately by GLB/non-GLB status across both the general population and GLB augmentations samples. For example, it can be applied when producing crosstabulations of a survey variable by the variable GLBstatus, which denotes whether or not a sample member is GLB-identified.

Weight 4: This weight should be used for analysis of all screened respondents separately by sample type, i.e., the general population sample or the GLB augmentation sample. For example, it can be applied when producing crosstabulations of a survey variable by the variable Recsource, which denotes the recruitment source for each case (with a value of 1 indicating that the case is from the general population sample and a combination of values 2-4 indicating that the case is from the GLB augmentation sample).

Weight 5: This weight should be used for analysis of all coupled respondents separately by sample type, i.e., the general population sample or the GLB augmentation sample. For example, it can be applied when producing crosstabulations of a survey variable by the variable Recsource, which denotes the recruitment source for each case (with a value of 1 indicating that the case is from the general population sample and a combination of values 2-4 indicating that the case is from the GLB augmentation sample).

Weight 6: This weight should be used for analysis of all screened GLB respondents who were active on the Knowledge Networks panel at the time of the survey and who were pre-identified as GLB prior to the survey, both those in couples and those not in couples.

Weight 7: This weight should be used for analysis of all coupled GLB respondents who were active on the Knowledge Networks panel at the time of the survey and who were pre-identified as GLB prior to the survey.

Supplemental notes on the weights from the Stanford Research Team:

Note 1: Couple Weights.

**NOTE ON COUPLES VERSUS INDIVIDUALS:** All of the weights **weight1- weight7** are weights based on the respondent only. In other words **weight1- weight7** are individual weights that count the respondents, but not the partners. If you want to count the partnered cases as couples, meaning two adults, you would have to divide the weighted count by 2 to get appropriate US national counts of couples.

The other option for couples is to use **weight\_couples\_coresident**, which started with a couple weight derived from **weight2**, and then reweighted the couples based on the cross classification of both partner's races, using the American Community Survey (ACS) of 2007 as the benchmark, and dividing couples into heterosexual married couples, heterosexual unmarried cohabiting couples, gay male cohabiting couples, and lesbian cohabiting couples. The supplementary **weight\_couples\_coresident** is only available for coresident couples, because the ACS only has information on both partners of a couple when the partners are coresident.

Stanford Research Team supplementary weight Note Two: Applicability of the weights to measurement of same-sex couples.

Because of the oversample of same-sex couples in the HCMST data, nearly 16% of the partnered adults in the HCMST have same-sex partners:

```
. tabulate same_sex_couple
```

best guess as to whether the couple is a same-sex couple	Freq.	Percent	Cum.
different sex couple	2,535	84.25	84.25
same-sex couple	474	15.75	100.00
Total	3,009	100.00	

Because GLB respondents were over-sampled in the design of the survey, GLB respondents are under-weighted in the weights to allow for nationally representative data to reflect the correct number and proportion of same-sex couples in the US. About 2% of all partnered adults in the US have same-sex partners.

```
. tabulate same_sex_couple [fweight=weight2]
```

best guess as to whether the couple is a same-sex couple	Freq.	Percent	Cum.
different sex couple	166,656,546	98.02	98.02
same-sex couple	3,360,860	1.98	100.00
Total	170,017,406	100.00	

\* If you want nationally representative estimates of the whole US adult population (with or without romantic partners), you must use the weights.

But what if you want the best nationally representative estimate for gay and lesbian adults only? Then the need to use the weights is not so clear. The weights in the HCMST survey were created by KN to bring the KN panel sample into line with the Current Population Survey (CPS) with respect to age, gender, race, home Internet access, and region. The KN weights do not reflect the particular demography of same-sex couples in the CPS. Furthermore, the CPS captures only coresident same-sex couples, and the number of identifiable same-sex couples in the CPS is modest. Lastly, researchers have come to recognize some of the inherent weaknesses and measurement error in the indirect way in which the US Census and the CPS identify same-sex couples.<sup>1</sup>

Furthermore, not all the same-sex couples in the HCMST data are from the over-sampled populations:

<sup>1</sup> See, for instance, Black, Dan, Gary Gates, Seth Sanders, and Lowell Taylor. 2007. "The Measurement of Same-Sex Unmarried Partner Couples in the 2000 U.S. Census." California Center for Population Research. <http://papers.ccpr.ucla.edu/papers/PWP-CCPR-2007-023/PWP-CCPR-2007-023.pdf>; and O'Connell, Martin, and Daphne Lofquist. 2009. "Counting Same-sex Couples: Official Estimates and Unofficial Guesses." U.S. Census Bureau. <http://www.census.gov/population/www/socdemo/files/counting-paper.pdf>; and O'Connell, Martin, and Gretchen Gooding. 2006. "The Use of First Names to Evaluate Reports of Gender and Its Effect on the Distribution of Married and Unmarried Couple Households." in *Population Association of America*. Los Angeles.

```
. table recsource same_sex_couple if qflag==1, contents (freq mean glbstatus mean weight2 ) row
col
```

recruitment source	best guess as to whether the couple is a same-sex couple		
	different sex couple	same-sex couple	Total
gen pop sample	2,334 .01114 70543.02871	43 .790698 21466.62791	2,377 .025242 69655.23517
glb augment sample	162 1 8144.240741	366 1 5426.756831	528 1 6260.530303
glb withdrawn sample	27 1 21257.74074	57 1 7606.263158	84 1 11994.2381
glb item refused sample	12 1 9649.25	8 1 2255.625	20 1 6691.8
Total	2,535 .089546 65742.22722	474 .981013 7090.421941	3,009 .229977 56502.95979

The numbers in the table are, from top to bottom in each cell, unweighted frequency, proportion of that unweighted frequency who reported being gay, lesbian, or bisexual, and the average of weight2. The “glb” categories are the recruitment sources that were oversampled, but 43 of the 474 same-sex couples come from the general population sample. When applying the weights, the 43 same-sex couples from the general population end up being over-weighted compared to the other same-sex couples, and this may not be the desired behavior. Having one subgroup of same-sex couples with much larger weights than the rest can skew estimates, reduce the effective sample size, and thereby reduce standard errors within the GLB or same-sex couple population. Given that the HCMST weights are not especially applicable to GLB adults and same-sex couples, we recommend:

\* When analyzing data for only the GLB respondents or only for same-sex couples, and if you want averages or coefficients for within-group analysis (rather than national totals) consider not using the weights at all.

What approach to weights is recommended for regression analysis comparisons of gay and straight respondents? Here the standard approach would be to simply apply the weights. A second approach follows Winship, Christopher, and Larry Radbill. 1994. "Sampling Weights and Regression Analysis." *Sociological Methods and Research* 23 (2):230-257, whose recommendation is that for many cases unweighted regression is superior to weighted regression as long as the regression predictor variables include the predictors of the weights.

\* In regression analysis comparing the heterosexual and the GLB respondents, using the weights can be appropriate, but it may also be appropriate to do unweighted regression, and include the predictors of weight in the regressions (following Winship and Radbill).

Predictors of the weights in HCMST include, most importantly: **recsource** followed by race, age, age squared, education, metro status, and Internet access at home (using the variable **ppnet**). See below for the regression results predicting **weight2**:

```
desmat: regress weight2 @ppage @age_sq ppmsacat ppnet recsource ppethm ppeducat
@pphouseholdsize
```

---

Linear regression

---

Dependent variable	weight2
Number of observations:	3009
F statistic:	112.829
Model degrees of freedom:	15
Residual degrees of freedom:	2993
R-squared:	0.361
Adjusted R-squared:	0.358
Root MSE	39372.129
Prob:	0.000

---

nr	Effect	Coeff	s.e.
1	respondent age at time of HCMST wave I survey	-723.353**	253.749
2	age_sq	5.173*	2.514
	ppmsacat		
3	metro	-6649.935**	2077.962
	ppnet		
4	yes	-7391.898**	1959.179
	Recsource (compared to general population sample)		
5	glb augment sample	-54486.781**	2143.554
6	glb withdrawn sample	-51213.847**	4437.261
7	glb item refused sample	-54409.463**	8900.359
	Ppethm (compared to white)		
8	black, non-hispanic	19951.284**	2794.533
9	other, non-hispanic	42840.931**	4108.487
10	hispanic	25571.866**	2415.068
11	2+ races, non-hispanic	-36984.609**	3807.468
	Ppeducat (compared to <HS)		
12	high school	-130.749	2676.084
13	some college	-11263.054**	2727.364
14	bachelor's degree or higher	-12863.442**	2668.883
15	number of people living in your HH	-895.552	537.434
16	_cons	106378.636**	7356.972

---

\* p < .05

\*\* p < .01