



Data Collection Worksheet

Please Note: The Data Collection Worksheet (DCW) is a tool to aid integration of a PhenX protocol into a study. The PhenX DCW is not designed to be a data collection instrument. Investigators will need to decide the best way to collect data for the PhenX protocol in their study. Variables captured in the DCW, along with variable names and unique PhenX variable identifiers, are included in the PhenX Data Dictionary (DD) files.

The Dissimilarity Index is based on U.S. Census Bureau data. This protocol describes how to make calculations using 5-year American Community Survey (ACS) estimates.

The ACS data used in this protocol can be accessed by using Excel to read the Summary Files or using the “Download Center” at the U.S. Census Bureau’s American FactFinder portal at <http://factfinder.census.gov>. Users can find additional information on these tools at the following locations:

Using Excel to Access Summary Files: http://www2.census.gov/programs-surveys/acs/summary_file/2014/documentation/tech_docs/ACS_SF_Excel_Import_Tool.pdf

Using the Download Center: http://www2.census.gov/programs-surveys/acs/summary_file/2014/documentation/tech_docs/How_to_Access_ACS_Estimates_AFF.pdf

The technical documentation for the ACS summary files is available online at <http://www.census.gov/programs-surveys/acs/technical-documentation.html>. Select the “Summary File Documentation” link, and then select the data set of interest. Users not familiar with Census data should consult the technical materials.

The key race/ethnicity data in the ACS are found in "Table B03002: Hispanic or Latino by Race." This table is preferred over other possible race and race/ethnic tables available, as it provides data on the main race/ethnic groups in the United States and explicitly incorporates data on Hispanic or Latino populations, otherwise not available in the race-only tables.

Variable Code	Variable Name
B03002001	Total:

B03002002	Not Hispanic or Latino:
B03002003	White alone
B03002004	Black or African American alone
B03002005	American Indian and Alaska Native alone
B03002006	Asian alone
B03002007	Native Hawaiian and Other Pacific Islander alone
B03002008	Some other race alone
B03002009	Two or more races:
B03002010	Two races including Some other race
B03002011	Two races excluding Some other race, and three or more races
B03002012	Hispanic or Latino:
B03002013	White alone
B03002014	Black or African American alone
B03002015	American Indian and Alaska Native alone
B03002016	Asian alone

B03002017	Native Hawaiian and Other Pacific Islander alone
B03002018	Some other race alone
B03002019	Two or more races:
B03002020	Two races including Some other race
B03002021	Two races excluding Some other race, and three or more races

The race/ethnic data are available for all small census geographies-such as census block, census block group, and census tract-and can be easily extracted for almost any geographic level. Note: Although block group data have long been available from the Census File Transfer Protocol site, the Census Bureau did not make block groups available for download at American FactFinder until the release of the 2009-2013 ACS. Information about accessing block group data for earlier years is available at http://www.census.gov/library/video/acs_block_group.html.

Researchers can use the data in this table to easily calculate basic variables (e.g., the percentage of any race and/or ethnicity group) or to combine groups (e.g., all minorities).

The Dissimilarity Index provides data on larger areas (e.g., metropolitan statistical areas) using smaller-level data.

The most common conceptualization of residential segregation is based on the dimension of evenness. *Evenness* refers to the differential distribution of the subject population across neighborhoods in a large area (e.g., metropolitan area). It ranges from 0 (complete integration) to 1 (complete segregation) and indicates the percentage of a group's population that would have to change residence for each neighborhood to have the same percentage of that group as the metropolitan area overall. It is computed as:

$$D = .5 * \sum_{i=1}^n \left| x_i / X - y_i / Y \right|$$

where

n is the number of tracts in the larger area (e.g., a metropolitan area),

x_i is the population size of the minority group of interest in tract i ,

X is the population of the minority group in the larger area (e.g., metropolitan area) as a whole,

y_i is the population of the reference group (usually non-Hispanic whites) in tract i , and

Y is the population of the reference group in the larger area (e.g., metropolitan area) as a whole.

The calculation requires the computation of the totals for each group across all subareas within a larger region (e.g., all census tracts within a county), the proportion of each group within each subarea, the absolute difference between the proportions, and the sum of the absolute differences. The latter number is multiplied by 0.5 to generate a result between 0.0 and 1.0. A value of 0.0 would indicate there were the same proportions of majority and minority group populations in each subarea, as in the larger regions' population. If all subareas within the region contain members of just one group (i.e., there is no co-residence), then D equals 1.0, indicating complete segregation.

Extending the Dissimilarity Index: The Multigroup Analog

While much early research on segregation looked at two groups (e.g., black and white, or majority and minority), today's society is multiethnic. Two-group measures are useful but limited for describing complex patterns of segregation. The choice to use a two-group or multigroup D depends on the specific question of interest. In a region where the population is composed of three groups (e.g., white non-Hispanic, black non-Hispanic, and Hispanic), we may be interested in

a) segregation between two specific groups (e.g., How segregated are white from black residents?); or

b) segregation among all three groups (e.g., How segregated are white non-Hispanic, black non-Hispanic, and Hispanic residents from each other?).

The two-group measure can still be used by comparing all possible pairs of population groups (Morrill, 1995), but these are not comprehensive, and multiple groups are not treated simultaneously. To address segregation among multiple groups requires a multigroup analog to D (Morgan et al., 1975; Sakoda, 1981). The multigroup analog describes the extent to which two or more population groups are similarly distributed among subareas. The formula for multigroup dissimilarity (from Reardon & Firebaugh, 2002) is:

$$D = \sum_{m=1}^M \sum_{j=1}^J \frac{t_j}{2TI} |\pi_{jm} - \pi_m|$$

where

T is total population,

M is the number of groups m ,

J is the number of subareas or units j ,

t_j is number of individuals in subarea j ,

π_m is the proportion in group m ,

π_{jm} is the proportion in group m , of those in unit j , and

I is the Simpson's Interaction Index, given by

$$I = \sum_{m=1}^M \pi_m (1 - \pi_m)$$

The interpretation of multigroup D (sometimes labeled as $D(m)$) is the same as D (Wong, 1993).

In the Stata statistical software package, the command `seg` (installed by typing "ssc install seg" from within Stata) will compute D (Reardon, 2002).

Researchers have extended segregation measures by incorporating the spatial dimension (White, 1983; Wong, 1993; Reardon & O'Sullivan, 2004). There are spatially modified versions of the D index (Wong, 1993).

Protocol source: <https://www.phenxtoolkit.org/protocols/view/211403>