## USP/PLSI 492: Research Methods

## Instructor: Pietro Calogero

# Getting Census data and prepping it for use in GIS:
# a real-world problem with Richmond, CA

Initially I wrote this tutorial using San Francisco as my example. However, that was misleadingly simple because the borders of U.S. Census tracts and county boundaries align. Therefore, with San Francisco—which is both a city and a county—you never have to deal with census tracts that overlap the border.

Richmond, CA, however, is only one city within Contra Costa County. And the borders of Richmond are disturbingly weird. The U.S. Census did not bother trying to get their tract-boundaries to align with the Richmond city borders at all. Census tracts are designed to have roughly the same population as each other, and designed to be fairly regularly-shaped. So in western Contra Costa County, they often overlap between city of Richmond, adjacent cities such as El Sobrante, and unincorporated county area.

Urban planners should be concerned about this because of the economic, racial, and health-inequities between Richmond and the wealthy communities on the west side of San Francisco Bay. Right away, we face much greater methodological challenges in studying Richmond, because erratic political boundaries make geopsatial analysis much more challenging. And that is why we should focus on Richmond, not already-gentrified enclaves like San Francisco.

## 1. Get the data.

Go to American FactFinder 2, Advanced tab, and set your filters to get data from the geographic area of your study, which whatever characteristics most closely correspond to your survey-design. Since Richmond does not correspond to census-tracts, and yet I want to work with tract-level data in Richmond, I go to the next-larger geographic unit: I get the Census Tract data for Contra Costa Co:



Select the table or tables that most closely correspond to the whole population (universe) that you are seeking to represent with survey data. As in this Richmond example, you may find that you need to go to a larger geographic unit, a larger population, and then refine the data by various criteria—including geography—until you can get a demographic profile of the population from which you want to do your sample-survey. In this

case I am interested in housing data, and I know that what I want is in Table DP04. Therefore, with just one filter—Census tracts of Contra Costa County—one of my first hits is the table I want. However, you may want to set many more filters so that you are not picking through 5,641 tables.



Once you have found a table that you think will have your data, click the hyperlink of the table name itself so that you can preview the data. Note: if you are getting ACS data, rather than straight-up decennial census data, the ACS data will be estimated values. And it will include margins of error, and percent margins of error for each data-point. That means that ACS data will show 4X the number of data-columns. To reduce this before downloading, click the "Modify Table" link in the Actions bar.



Then, in the top left cell of the table itself (above the word **Subject**), a small blue funnel will appear. Click on this funnel, and check only the "Estimates" column. This will hide the "Estimate Margin of Error", "Percent", and "Percent Margin of Error" columns, and exclude them from the data-download. That reduces the amount of data you will need to sift through by 75%.
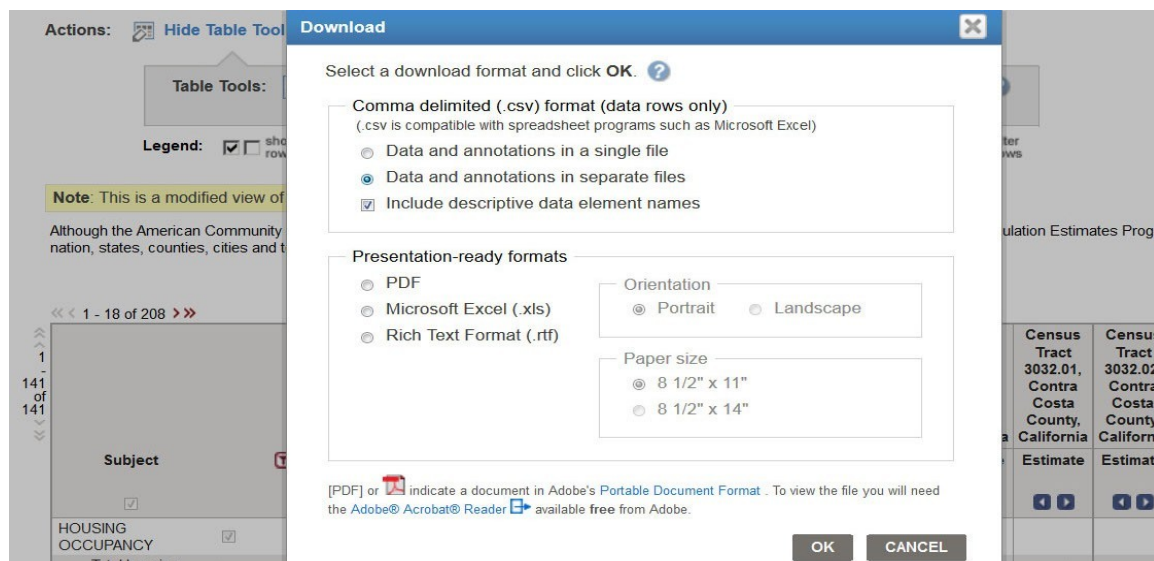
Note that before I reduce the desplay to "HC01/Estimate" only, there are 832 columns in this table...
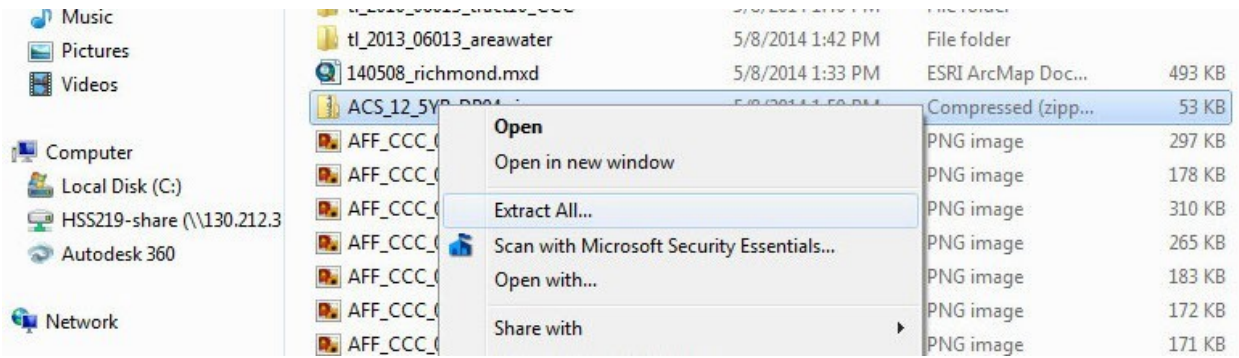


Once I filter the view to "Estimates" only, the number of columns drops to 208.

Select "Download…" And then choose the .CSV option. That will be a text file with values separated by commas. American FactFinder also offers an option to download the file in "Presentation-ready formats". Do not succumb to the temptation of convenience! If you pick any of these options, it will strip off geocodes and it will be far more difficult to join this tabular data with TIGER geospatial data.



DO choose "separate files", and DO choose to "include descriptive data".

Save the file at a PATH/NAME that you will be able to find again. And remember to extract it from the .zip archive:
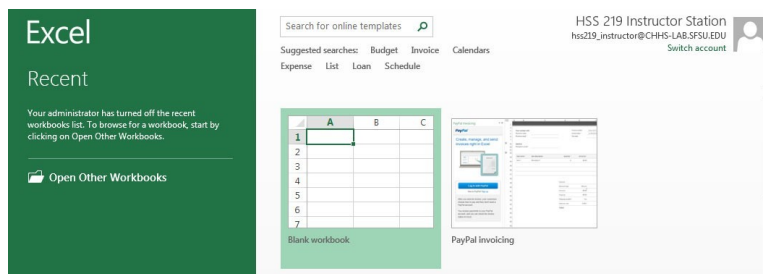


## 2. Prepare and arrange the data in Excel.

ArcMAP and QGIS can import .csv files. But the file you download from FactFinder will need to be cleaned up in two ways, which I call "winnowing" and "relabeling". Before we winnow and relabel, though, make sure you learn how to bring CSV data into Excel in a way that does not corrupt the data.
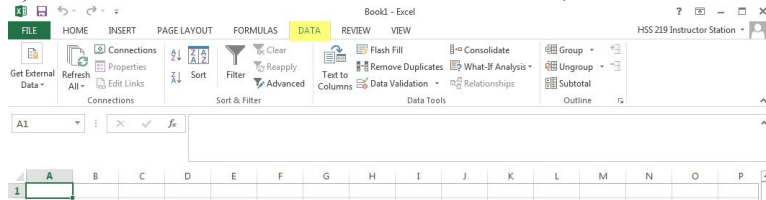
### 2a. Import a CSV into Excel without losing leading zeroes.

Importing CSV data into Excel is also a little tricky. If you open a FactFinder CSV file directly into Excel, it will strip off any leading zeros in any data. In California, the GEO.ID for census tracts begins with a zero, such as: 06075010100.
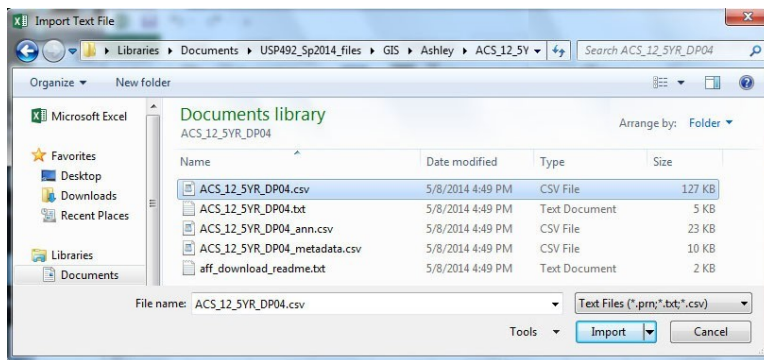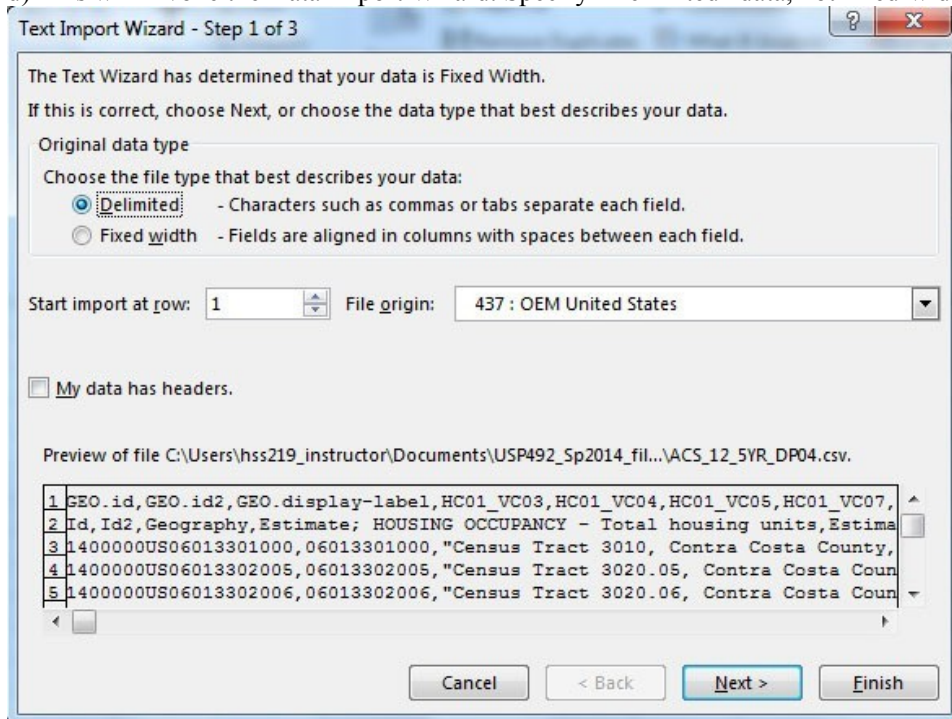
So:    a) Open a **blank** Excel spreadsheet.



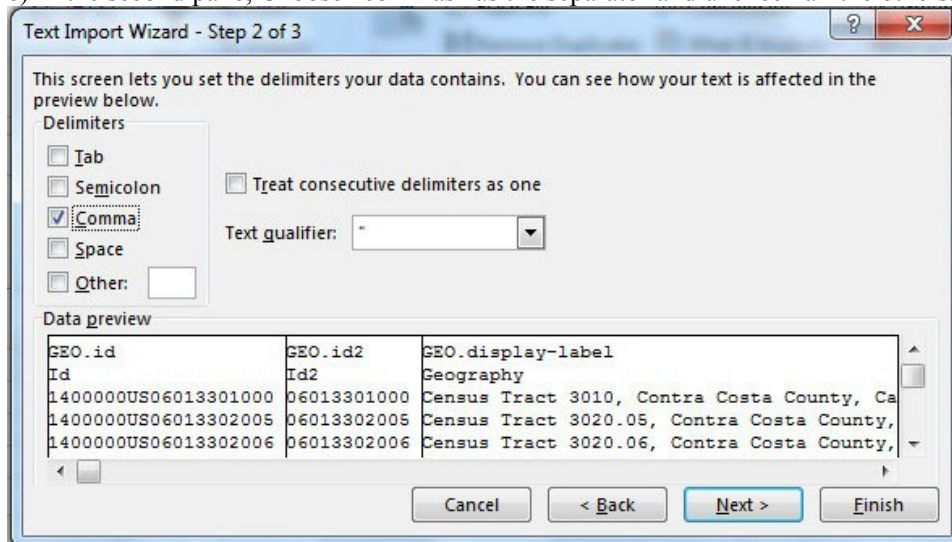b) Under the Data tab > Get External Data subtab, Select "From text."



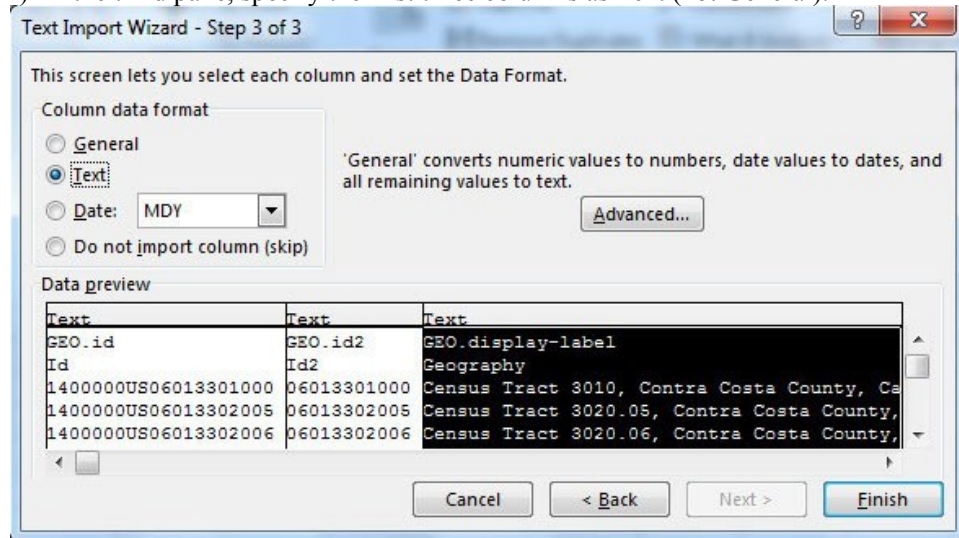c) Choose the Path and file. The largest CSV is usually the right one.

d) This will invoke the Data Import Wizard. Specify "Delimited" data, not fixed-width.
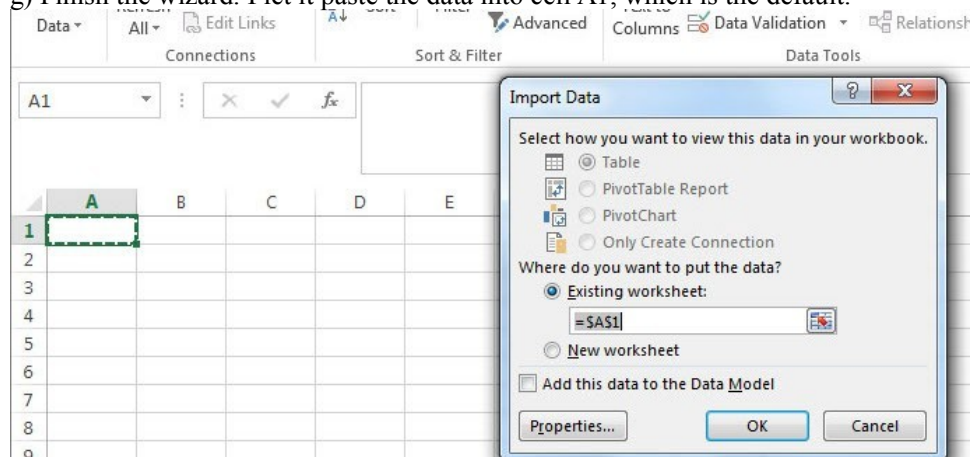


e) In the second pane, Choose "commas" as the separator and uncheck all the others.

f) In the third pane, specify the first three columns as Text (not General).



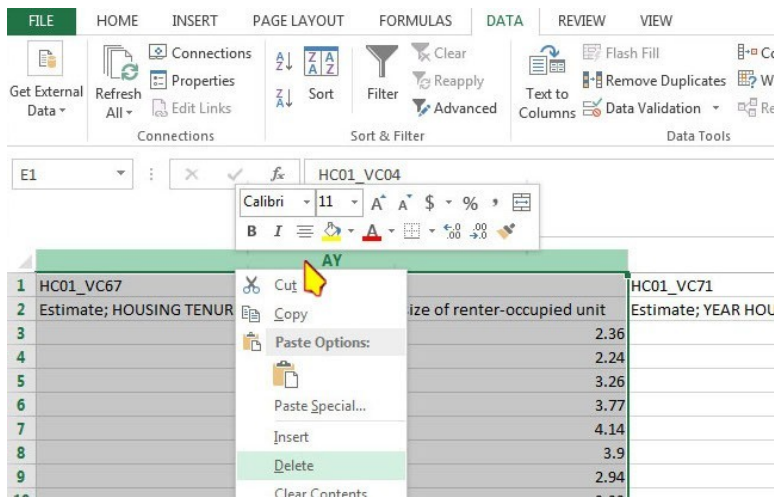g) Finish the wizard. I let it paste the data into cell A1, which is the default.



Check that the data imported cleanly.

## 2b. Winnowing the data: pulling the needle out of the haystack.
Table DP04 has more than 100 columns of data, covering all sorts of housing characteristics. I only want to analyze a few characteristics, so I delete most of the columns. What I keep are:

- Column B: GEO.id2, which is the "correlation column" is will use to join this data to a TIGER file.

- Column C: this names the Census tract. I keep this as a verification-check that the Table-Join works correctly.

- Columns D: the total number of housing units in each Census Tract.

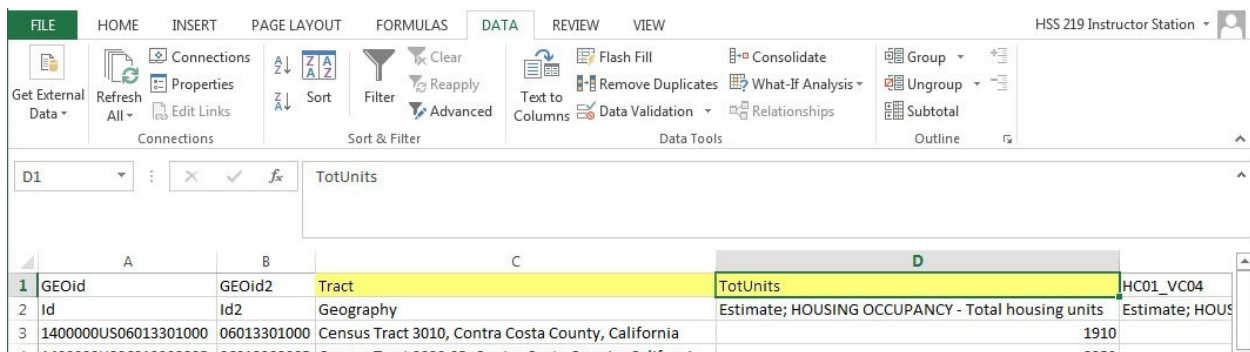- Any remaining columns that include data I will analyze.

One more thing: Excel creates three-sheet Workbooks by default. Just to reduce confusion later, go to the bottom left of the window-border select Sheet2, R-click, Delete; and Sheet3, R-click, Delete. You can eve rename the Sheet1 as "Data" or something to remind you of which sheet to look for as you import this into other programs.

## 2c. Re-labeling the headers

In this third Excel step we need to do two things: edit the header-row (Row 1) so that ArcMAP can read the data, and also edit the header labels so that we ourselves can understand what data we are looking at.

To make the headers readable by ArcMAP, any punctuation-marks must be removed (except under_score). Rename "GEO.id2" to "GEOid2". Rename "GEO.display-label" to "tract_num".



To make the remaining headers human-readable, rename them using the descriptive information in Row 2. [Back when you were downloading this CSV from FactFinder, the Row 2 information was included when you check the box "Include descriptive data element names".]
Replace the Row 1 codes with "dehydrated" but meaningful terms like "totalunits", "medvalue10", or "mvin2012". Max length: 12 characters. Don't start with a number. Don't use punctuation marks other than underscore: "_". Don't use spaces. Try to create the shortest possible name that you can intuitively understand.

Finally, I delete Row 2. Now, this table only has one header row (Row 1) and the rest of the rows are data-records for each census tract.

## 2d.  Use find/replace to shorten descriptive data.

The data-column that includes the tract-number in FactFinder files has a lot of text. Use the "Find/Replace" function to shorten the name:

I left the "Replace with" field blank, so that Excel deleted all of:
 , Contra Costa County, California



## 2e. Save as both an .XLS and a .CSV format

Just to cover your bases, save in both formats. ArcMAP can read Excel spreadsheets directly. QGIS cannot.

Here is another tip: Save the file with a very short filename. When you do a table join in GIS, it will concatenate the filename and the column-name together. So choose the shortest possible filename that still makes sense to you. In this case, "DP04.csv" would be a good name, because it reminds me where this data came from.