## TUTORIAL 6 | COMPARING CENSUS VARIABLES

**Goals**

- Become familiar with ACS data.
- Learn how to derive an average from census data.
- Consider the strengths and weaknesses of census data.
- Compare two different census categories geographically.
- Analyze demographic patterns through maps.

**Introduction**

In this tutorial, we'll be looking at a different kind of census data collected by a survey called the ACS (American Community Survey). After the 2000 census, the ACS began to replace the long-form decennial census. Now, every 10 years we answer only the short-form census; to supplement that data the ACS polls about 1% of the US population each year (around 3 million people). The ACS asks questions especially interesting to architects, such as age and type of housing, rental costs, occupants per room, mortgage costs, HVAC characteristics, and even building material. You can see the information ACS collects here:
https://www.census.gov/programs-surveys/acs/guidance/handbooks/general.html

Table excerpt from Chapter 1 below:



Table 1.1. **Population and Housing Data Included in American Community Survey Data Products**

| Social Characteristics | Economic Characteristics | Plumbing Facilities[6] |
|---|---|---|
| Ancestry | Class of Worker | Rent |
| Citizenship Status | Commuting (Journey to Work) | Rooms/Bedrooms |
| Citizen Voting-Age Population | Employment Status | Selected Monthly Owner Costs |
| Disability Status[1] | Food Stamps/Supplemental Nutrition Assistance Program (SNAP)[4] | Telephone Service Available |
| Educational Attainment | | Tenure (Owner/Renter) |
| Fertility | Health Insurance Coverage[2] | Units in Structure |
| Grandparents as Caregivers | Income and Earnings | Value of Home |
| Language Spoken at Home | Industry and Occupation | Vehicles Available |
| Marital History[2] | Place of Work | Year Householder Moved Into Unit |
| Marital Status | Poverty Status | Year Structure Built |
| Migration/Residence 1 Year Ago | Work Status Last Year | |
| Period of Military Service | | **Demographic Characteristics** |
| Place of Birth | **Housing Characteristics** | Age and Sex |
| School Enrollment | Computer and Internet Use[5] | Group Quarters Population |
| Undergraduate Field of Degree[3] | House Heating Fuel | Hispanic or Latino Origin |
| Veteran Status[2] | Kitchen Facilities | Race |
| Year of Entry | Occupancy/Vacancy Status | Relationship to Householder |
| | Occupants Per Room | Total Population |

[1] Questions on Disability Status were significantly revised in the 2008 survey to cause a break in series.
[2] Marital History, Veterans' Service-Connected Disability Status and Ratings, and Health Insurance Coverage were added in the 2008 survey.
[3] Undergraduate Field of Degree was added in the 2009 survey.
[4] Food Stamp Benefit amount was removed in 2008.
[5] Computer and Internet Use was added to the 2013 survey.
[6] One of the components of Plumbing Facilities, flush toilet, and Business or Medical Office on Property questions were removed in 2016.
Source: U.S. Census Bureau.

In this tutorial, we'll be looking at two variables: Home Occupancy Status (B25003), which specifies "renter" or "owner", and Median Household Income (B19013) over the last 12 months.
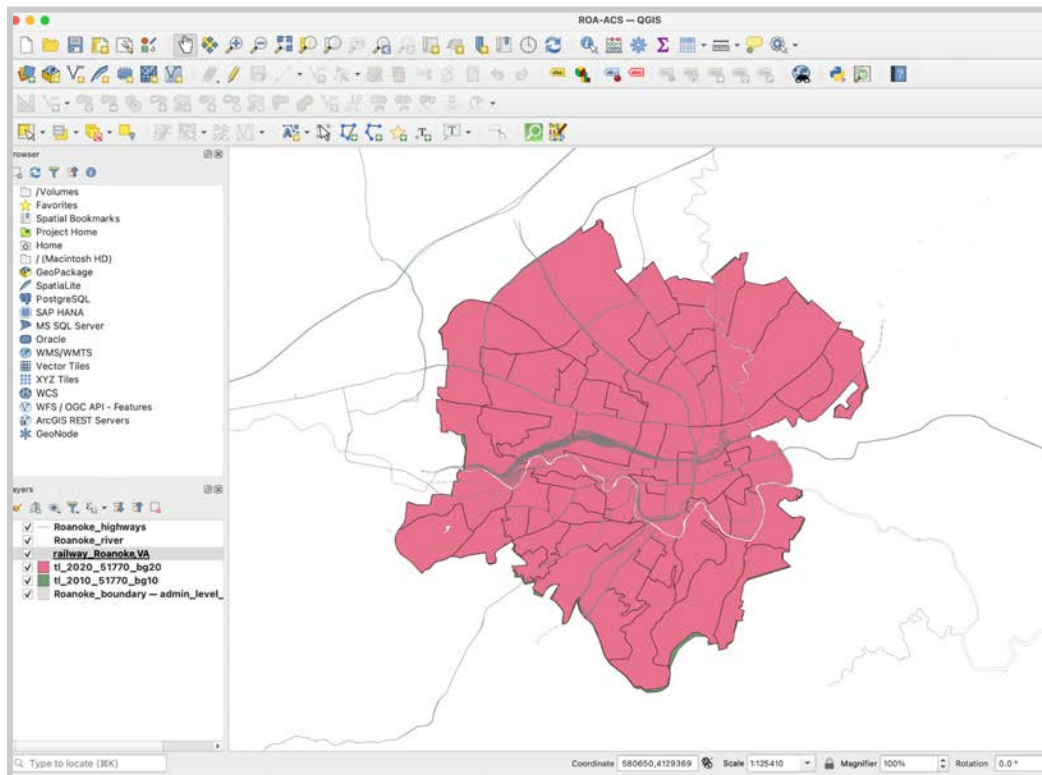
ACS data are aggregated to adjust for their small sample size. They come in 1-year, 3-year, and 5-year estimates. For these tutorials, we'll be using the 5-year estimates, which combine five adjacent ACS surveys. For example, the 5-year estimate for 2019 uses data collected from 2015-2019 – and not only 2019. These range estimates make ACS fundamentally different from decennial data which represents a single snapshot in time. Keep this in mind when interpreting and labeling the data.

**Step 1: Open a new file and import TIGER/Line shapefiles and physical geography layers.**

1a Set your CRS to 17N.

1b Import your city Block Group TIGER/Line Shapefiles for 2010 and 2020.

1c Add a few physical geography layers for reference. I recommend adding the railroad lines and major roads / highways, as well as the river.
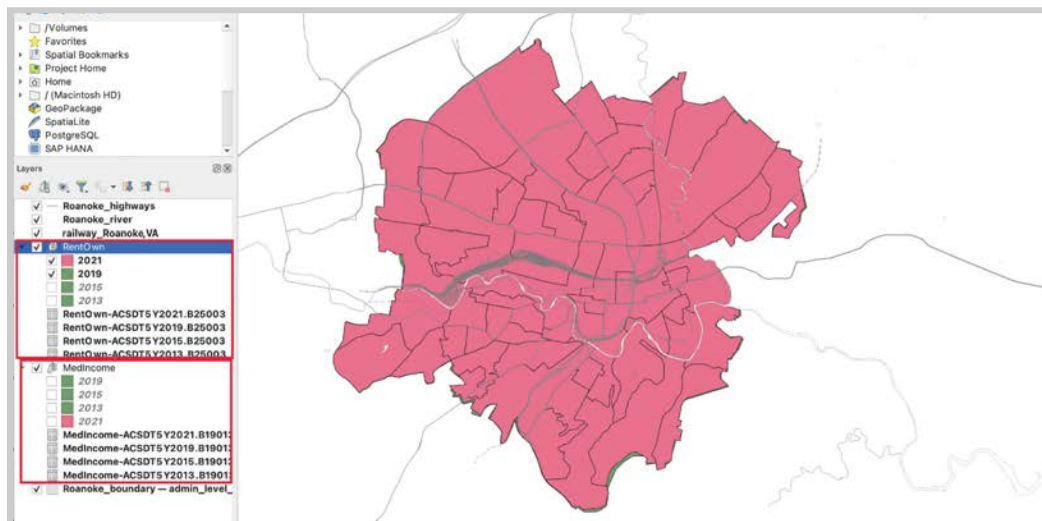
**Step 2: Associate new ACS data with the TIGER shapefiles.**

You'll see that you have 4 ACS datasets for each variable: the 5-year estimates for 2013 (the earliest year available), 2015, 2019, and 2021 (the most recent available). This range will allow you to see patterns over time, though keep in mind that each estimate includes the previous 4 years as well. So, the 2013 dataset includes information collected between 2009-2013; the 2021 includes 2017-2021; etc.

The **2013, 2015, and 2019** estimates will all join to the **2010 shapefile**. Any estimate between 2010-2019 will use the 2010 shapefile. Estimates beginning in 2020 use the 2020 shapefile, so the 2021 data will join with the 2020 shapefile. You'll make a separate copy of your shapefile layer for each dataset joined, rather than joining multiple datasets to one shapefile. This will make it easier to turn the different years on and off.

2a Import the **4 data tables each** for B25003 (Rent vs. Own) and B19013 (Median Income) (**8 csv files total**).

2b Create 5 copies (right click > **duplicate layer**) of your 2010 block group shapefile and 1 copy of your 2020, so that you have **8 shapefile layers** total to match your 8 data tables – 6 from 2010, and 2 from 2020. **Group** half the shapefiles with one ACS variable and half with the other. You should end up with **one group for Median Income** (4 data tables and 4 shapefiles) and **one group for Homeownership Status** (4 data tables and 4 shapefiles).

Now, **re-name your shapefile layers** within each group to match the year of the data table you'll join with it (eg, "2013", "2021"). This will help you keep track of the joins.



2c Join your .csv layer with your shapefiles (**double click > Join; GEOID to GEOID**). Make sure to match each year to the appropriate layer. Make sure to check that the join worked by looking at each layer's Attribute Table afterwards and making sure you don't see a column of "Null" values. Make sure to join the 2021 data tables to the 2020 shapefiles, and all others to the 2010 shapefiles.

BE CAREFUL TO JOIN THE
CORRECT YEAR AND VARIABLE
WITH THE CORRECT SHAPEFILE

**Step 3: Style the data according to median income.**

Like we did in the last tutorial, we need to style the two variables in this tutorial – median income and homeownership status – using a consistent style so we can compare between years. To take the example of median income, the range of median incomes varies between years. In 2013 in Roanoke, the median income for a given block group ranged from $6k - 147k, while in 2021 it ranged from $20k - $146k. That's a pretty large increase in the lowest median income. Keep in mind that these are median household incomes, not median individual incomes.

There are many things that might account for this increase. It could be that the census changed their estimation method, or that the lower-income areas grew in population and raised the number. Typically census block groups are organized roughly around population, so the population should be fairly consistent from one block group to the next. However, it's worth considering that a very high or very low income block group – especially one that varies considerably from the next closest median income – might have an especially low population that skews the data. You can easily see this by first checking your Attribute Tables and organizing them by median income to check for outliers (unusually high or low values). Second, you can check the population of each block group using the "Total" column in the *Homeownership Status* data.

3a To start off, let's find an appropriate range for the median income data. First, open each of your **four median income data table csvs** in Excel, Numbers, Google Sheets, or another spreadsheet software. Order each by its **median income column**, both **Ascending and Descending**. Make a note of the **highest and lowest values of each year** (2013, 2015, 2019, and 2021), and then take the overall lowest and overall highest number as your **min and max**. For Roanoke, I found a min of $6029 (round **down** to $6,000) and a max of $155208 (round **up** to $160,000) – so my range is 6k - 160k.
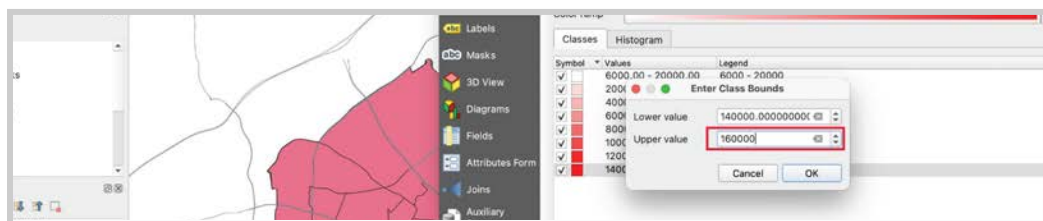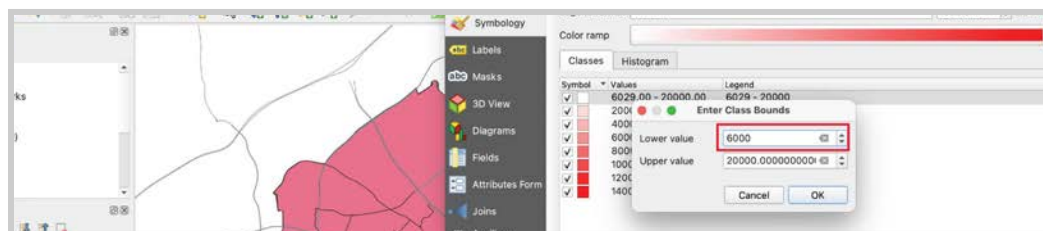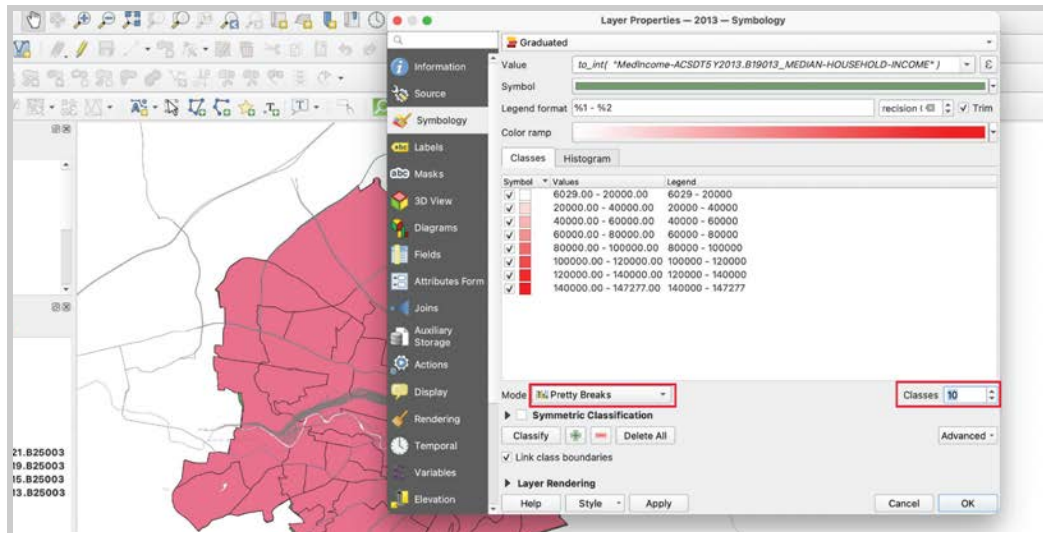
3b Next, open the **Symbology** in your Layer Properties for the median income shapefile with the **lowest minimum**. Select "**Graduated**", and then open your **field calculator** for Value. Expand your "**Fields and Values**" and find the Median Income variable. Notice that beside the variable is a small "**abc**", rather than a "123". This means that the variable is a "string", or text-based, rather than an "integer", or number-based.

You'll need to convert the variable to an integer before you can use it with the graduated style. To do this, expand the "**Conversions**" section of the field calculator and double click on "**to_int**" near the bottom. This will insert a "to_int ( " to your field calculator. Now, add in the **Median Income** field (double click it under Fields and Values), and the closing parentheses – it should be **to_int ("MedianIncome")**
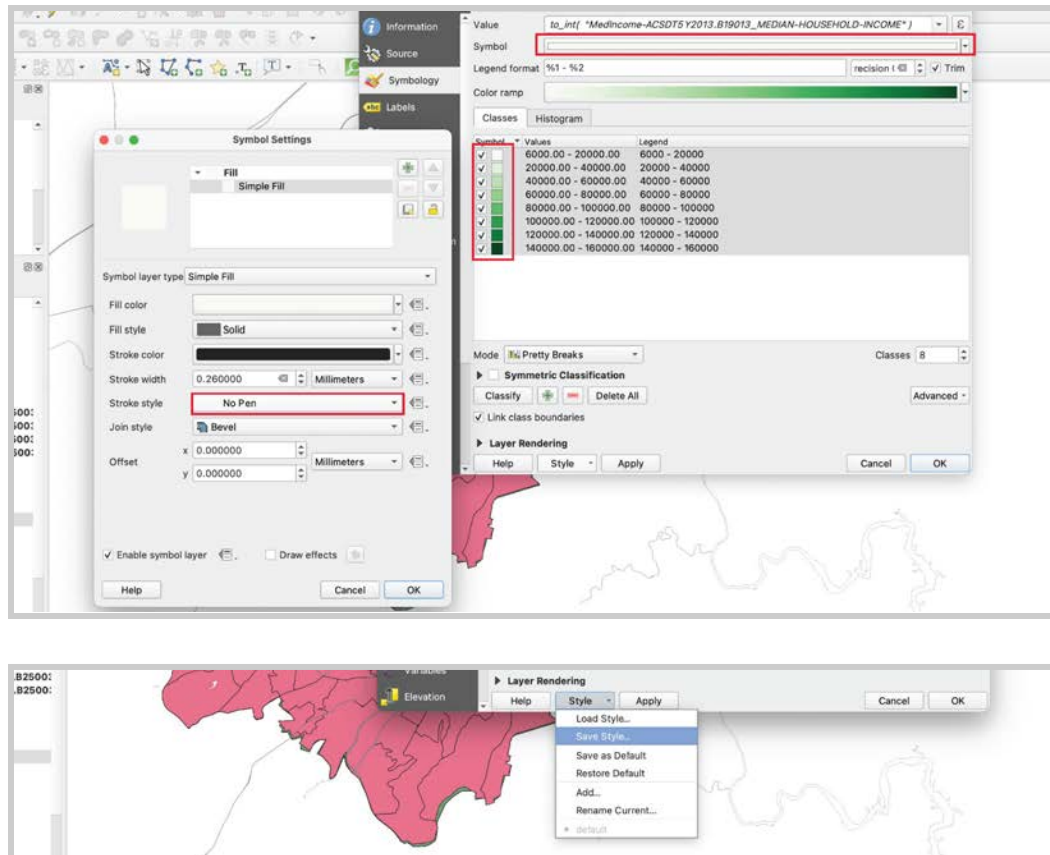


As always, check that the "**Preview**" looks alright, and then click "Ok".

3c This time we're going to use the classification **mode** "**Pretty Breaks**" to set up our intervals. First set your classes to 10, and then select "Pretty Breaks". You'll see that it produces pleasingly even numbers. After this, **manually set** your **minimum** value to the lowest median income you found across *all four data sets,* and your **maximum** value to the highest you found.
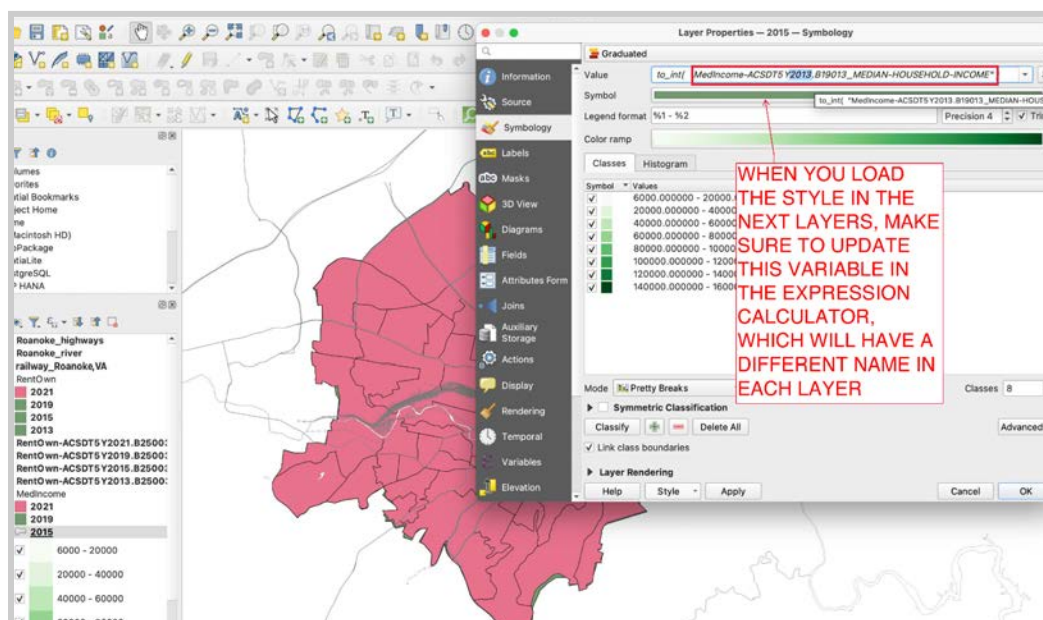
Select all classes and set your boundary style to **"No Pen"**, and then **Save your Style**.

NOTE: Make sure that each variable has its own color gradient. Use a different gradient for race vs. income vs. ownership status. This is a cardinal rule of mapping: **use a different color for each variable**.

3d Apply the style to your other three median income layers (**Load Style**). Note that you will need to **change the field name in the expression calculator** to match each data table's differently named variables. Usually this will only require changing the year (eg from "2013" to "2015") in the variable equation.

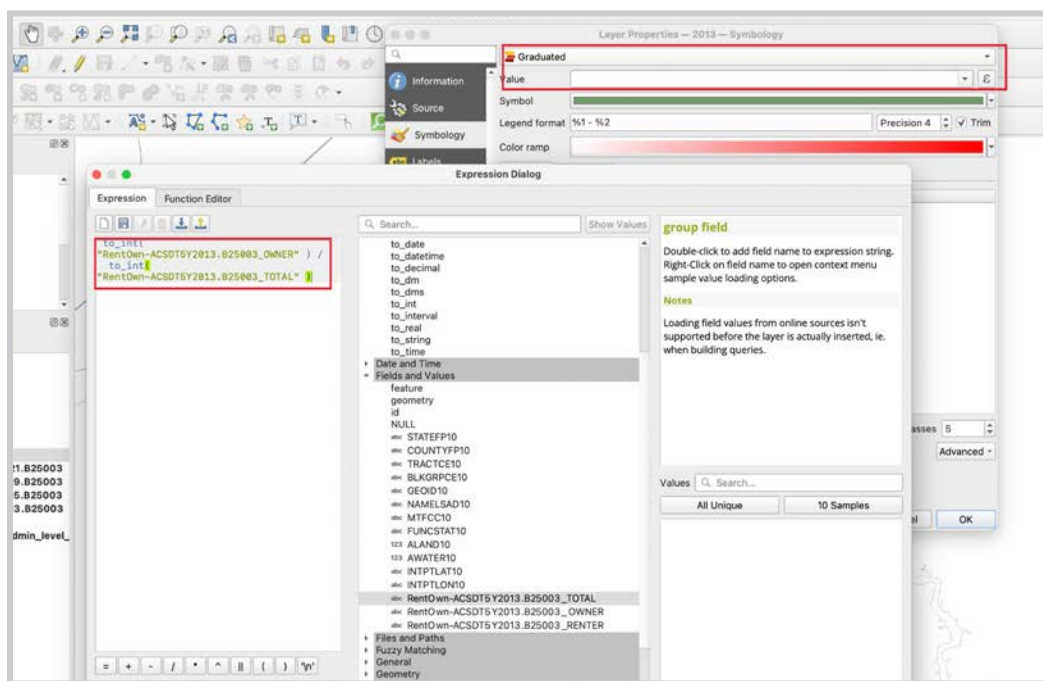**Step 4: Join and style your second ACS variable, Homeownership Status.**

4a Join your B1903 Homeownership data tables with their corresponding shapefiles, as you did for median income. Check the shapefile Attribute Tables to make sure they joined correctly.

4b You should now have two layer groups: one for Median Income and one for Homeownership Status, each with four joined shapefile layers: 2013, 2015, 2019, and 2021.
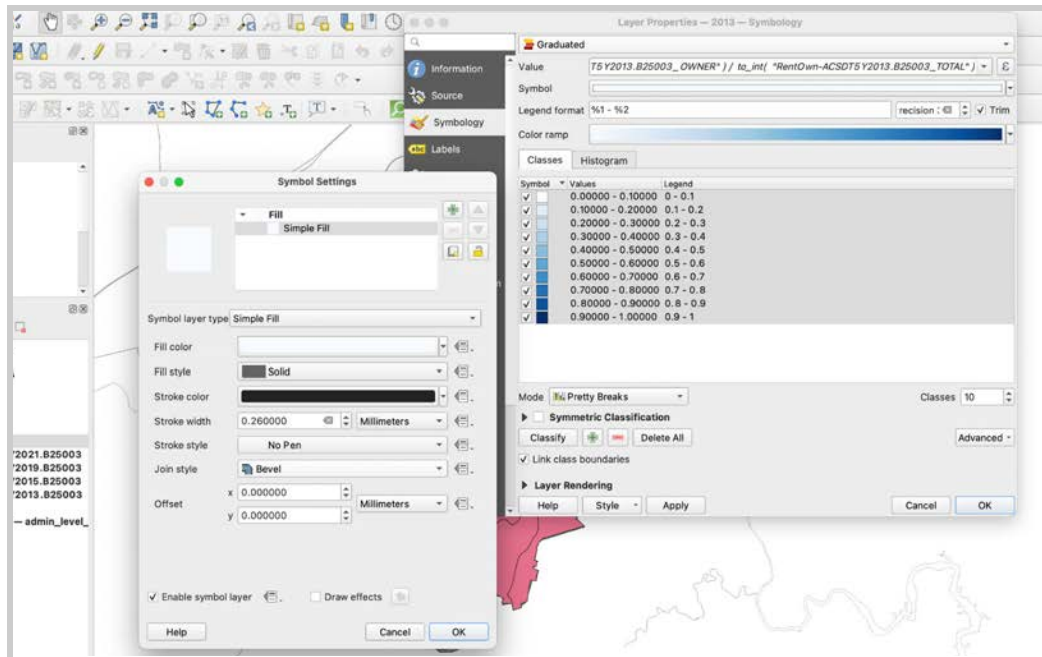
Now, open up one of your Homeownership shapefiles and set the Symbology to "**Graduated**". As with Median Income, you'll need to use the field calculator for this one. Open it up and check your "**Fields and Values**". Notice that, again, the Homeownership fields (Total, Owner, and Renter) are "abc" text fields rather than integer fields. This is a common issue with ACS data. Luckily, the solution is simple: we can use the to_int formula.

For the Homeownership data, like for the Race data in the previous tutorial, we want to look at percentages rather than raw numbers. This will help us adjust for population differences between the block groups. To get the **percentage** of home ownership, we'll **divide the "Owner" variable by the "Total" variable**. You'll need to convert them both to integers first (**Conversions > to_int**). So, your formula will look like: to_int ("Owner") / to_int ("Total"), where "Owner" and "Total" are the names of your specific fields.
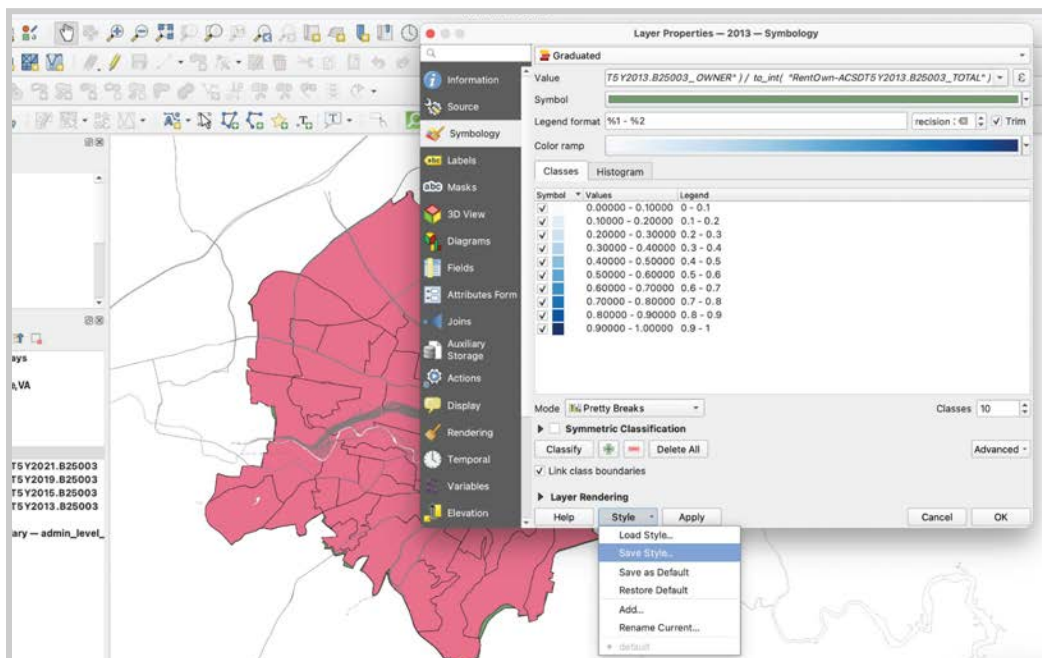
Before clicking OK, check that the "**Preview**" looks alright.
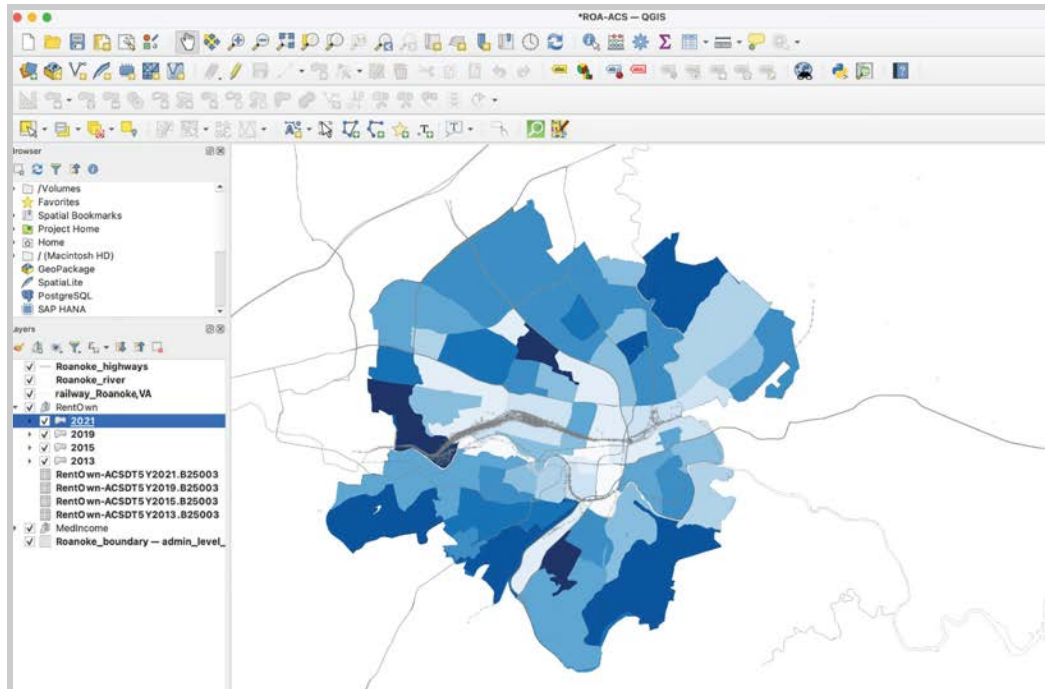


4c Set your classes to "**10**", and select "**Pretty Breaks**" again. This may get you even 10 percentage point breaks (0-0.1, 0.1-0.2, etc.), but you may need to **manually** edit some of the values. In my case, the lowest number was 0.04 instead of 0, so I manually changed that value to 0. Make sure you have 10 classes, each with a 0.1 interval, from 0 to 1 (0% to 100%).

4d Make sure to set your outline style to "**No Pen**" and choose a different color gradient than you used for Median Income. **Save your Style.**



4e **Load your Style** for each of your three remaining Homeownership layers, being sure to update the Fields in your expression calculator to match.

**Step 5: Create a Print Layout to compare the two data sets.**

Flipping between the four data layers for each variable, you might see some slight trends upwards or downwards, but specific block groups (roughly, neighborhoods) might go up while the city goes down, or vice-versa. You may not see clear or explicit trends in this data.

However, if you start to compare between the two variables, you might notice more of a pattern. Unsurprisingly, lower median income corresponds with lower rates of home ownership. You can demonstrate this in your print layout by comparing your 8 maps, divided by year.

You might also notice one or two neighborhoods where the two variables do not correlate across all four years. Call one of these out in zoomed-in maps on a **second** print layout page. Using Google, identify the neighborhood and give a brief historical explanation of why median income and homeownership might not correlate.

Remember to add titles, subtitles, legends, north arrows, scales, bylines, and sources to each of your print layout pages.

**-    Bonus    -**

**Step 6:** Add a page to your layout to compare one of the ACS variables in this tutorial with the P1 decennial race data from the previous tutorial.

6a Copy your 2010 and 2020 block group shapefile layers and import the data tables you used in the previous tutorial. Join them to their respective shapefiles. Load the Style you saved in the previous tutorial for each layer, being sure to update the changed Field name in the field calculator.

6b Compare the 2010 P1 data with the 2013 ACS data, and the 2020 P1 data with the 2021 ACS data.

Outline specific zones to compare between the two data sets. What correlations do you see between race and the ACS variable?