## TUTORIAL 5 | USING CENSUS DATA

**Goals**
- Parse regional and local census data.
- Link census data to TIGER shapefile.
- Examine opportunities and limitations of census data.
- Compare variables across time.
- Produce choropleth map of chosen census feature.
- Add multiple maps at different scales to a print layout.

**Introduction**

In this tutorial, you'll be looking at State and City census data. You will look at the decennial census's racial information, and compare that data 1) between geographies (contextualize your city within the state), and 2) between years (contextualize today). In this tutorial we'll only use one census question: P1, the race question. We will not also be looking at H1, the ethnicity question (hispanic or non-hispanic), but keep in mind that the P1 question does not consider ethnicity, so it presents a limited picture of diversity.

As you've learned GIS data typically has two components: the data itself, usually in spreadsheet form, and the geographical shape that the data links to. The shapefiles are vectors, which means they are defined as points connected by lines (eg, the outline of a census district or building). Each polygon within a shapefile has an associated ID number. The data spreadsheet will have a corresponding ID column that links data (aka "attributes") with shape.

To use the Census data, we need to download both the shapefile (known as a TIGER/Line Shapefile. TIGER stands for Topologically Integrated Geographic Encoding and Referencing), and the spreadsheet of census information. The spreadsheet will link to the shapefile by associating their shared ID numbers. In this case, each polygon of a shapefile – whether a state, county, census district, and so on – has an official ID number. Typically, the counties within the state begin with the state's ID number, and the census tracts within the county begin with the county ID number, and so on. For instance, Virginia's ID number is 51, and Roanoke County's ID number is 51770 (51 + 770), and Roanoke's 10th Census Tract has ID number 517700010. Note that, because the ID numbers correspond to a certain size of data, you must also use the corresponding TIGER file. So, you must link county-level census data to the county-level TIGER shapefile.

The order of size is:
country > state > county > tract > block group > block

Census data is measured from the state level all the way down to the block level, which generally correspond to single city blocks. Not all data is available at every level. More importantly, the smallest level of data is not always the best for a given purpose. If you use block-level data when looking at a variable across an entire county, for instance, you won't be able to see meaningful patterns because the geographic unit is too small. Typically, you only want to use 1 or maybe 2 levels smaller of a dataset to look at patterns. So for instance if you are looking at data across the United States, you might look at the state or census tract level.

**Note 1: your census data and shapefile year must match.** The boundaries of geographical areas (census tracts, blocks, etc) CHANGE every 10-year census (so, your 2020 census data won't accurately match your 2010 shapefile). We'll be using both the 2010 and 2020 shapefiles in this tutorial to compare changes across the decades.

**Note 2:** the US government does a major, Decennial census every 10 years. The American Community Survey (ASC) supplements the major census by collecting annual data through a smaller, randomized poll that is then aggregated and used to predict the un-polled population changes. We'll use Decennial data in this tutorial, and talk more about the ACS in future tutorials.

**Step 1: Download US Census shapefiles and data from the Data folder.**

1a In the Data folder, you'll find the 2010 and 2020 TIGER/Line shapefiles for the state of Virginia, as well as for your individual city. Download both the "Virginia Census" folder, as well as the census folder for your city. As always, unzip and keep all files in all folders, both primary and helper files.

The state data we'll be looking at is broken down to the county level, and the city data is at the block group level. You'll see once you bring the data into QGIS that these scales make it easier to see patterns in the data. The folders you download contain both the shapefiles and the spreadsheet data, which you'll join in QGIS using their shared GeoID fields.

The Shapefile Data was downloaded from
https://www.census.gov/geographies/mapping-files/time-series/geo/tiger-line-file.html
and the csv files downloaded from
https://data.census.gov/

For more information on how to download and clean census data, which can be an involved process, see the later tutorial on this subject.

**Step 2: Import the State data to QGIS.**
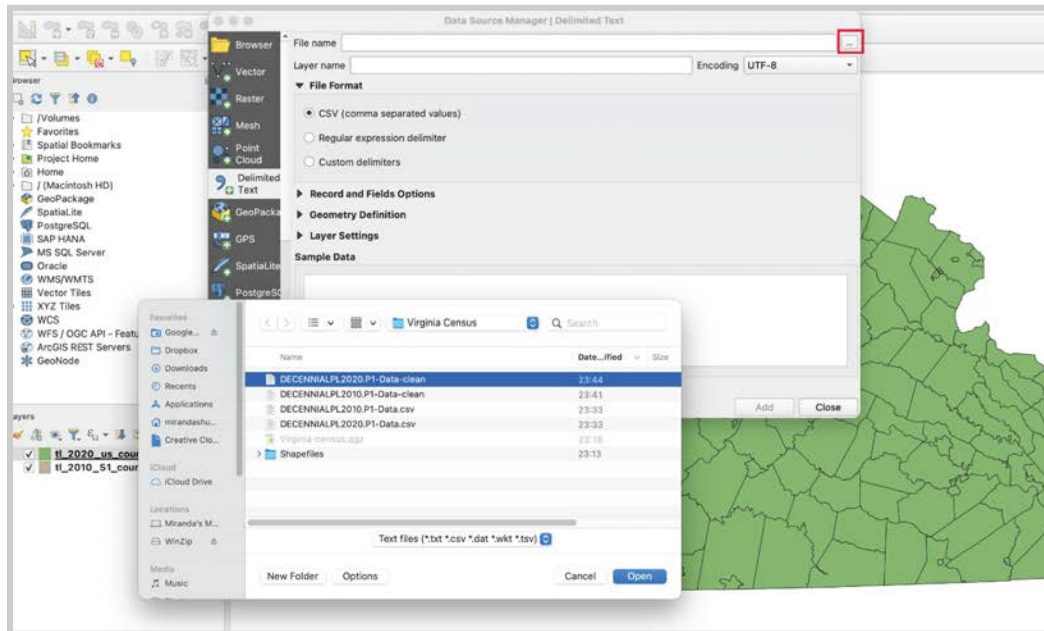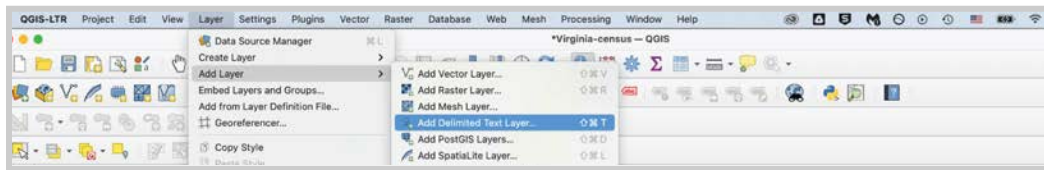
2a Open a new QGIS file.

2b Set your **CRS** in Project > Properties to **NAD / UTM 17N**.



2c First, add the state of Virginia shapefiles. Go to **Layer > Add Layer > Add Vector Layer**… and select the 2010 shapefile. Click "Add". Then, the 2020 shapefile and "Add"



2d Next, add the spreadsheet data, the .csv files, for the state. Go to **Layer > Add Layer > Add Delimited Text Layer** and select your **2010 and 2020** P1 csv files. Make sure that under "Geometry Definition" that **"No Geometry"** is selected. Check that the Sample Data display looks ok and click **"Add"**.

**Step 3: Associate Census data with TIGER shapefiles.**

3a First, we're going to open the Attribute Tables of the shapefile layers and the data table (csv) layers and compare them. Open the Attribute Tables for the 2010 shapefile and the 2010 csv.

Look over the fields in both files. You'll see that the shapefile has some standard information on county name, number, and area. The csv file has a few fields on racial information. In this tutorial we'll only use the "White" and "Total" columns, but you might look at other racial breakdowns for your own final project. This csv file is significantly reduced from the file you download from the census factfinder website, which contains well over a hundred fields breaking down race into combinations of 1, 2, 3, or more racial identities, plus margins of error.

Notice that the csv also has an ID field: GEO_ID. If you align the two Attribute Tables side by side, you'll see that this column matches the GEOID field in the shapefile. You can order both Attribute Tables by this column (click on the column name to sort the table) to check that the GEOIDs match. You will use this GEOID field to match the rows of data in the csv with the rows of data in the shapefile, thereby expanding the data associated with the shapefile geometry – along with area and name, this supplemental table will also add race information to the county shapes. Check that the GEOID column in the Attribute Tables in the 2020 files match as well.



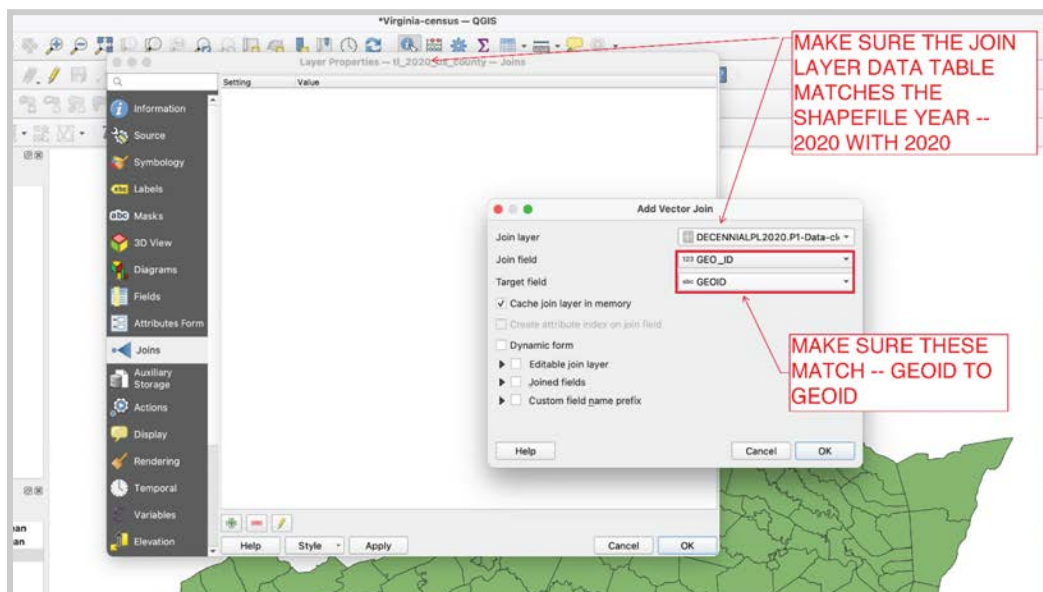> NOTE: The reason you can't join both 2010 and 2020 csvs to the same 2020 or 2010 shapefile is because the geographic boundaries of the census tracts, counties, blocks, and so on change slightly every decade. If you look closely, you'll see differences between the 2010 and 2020 county shapefiles.

3b To join the csv with the shapefile, open the **2010 shapefile** Layer Properties (double click) and select the **"Joins"** tab (blue arrow). Click the "Add" symbol at the bottom.

Select the 2010 csv as your "**Join Layer**". Then make sure that the "**Join Field" and "Target Field"** are both set to **"GEOID".**



3c Open the 2010 Shapefile Attribute Table, and make sure that you don't see "NULL" for the fields (this happens when the merge fails).
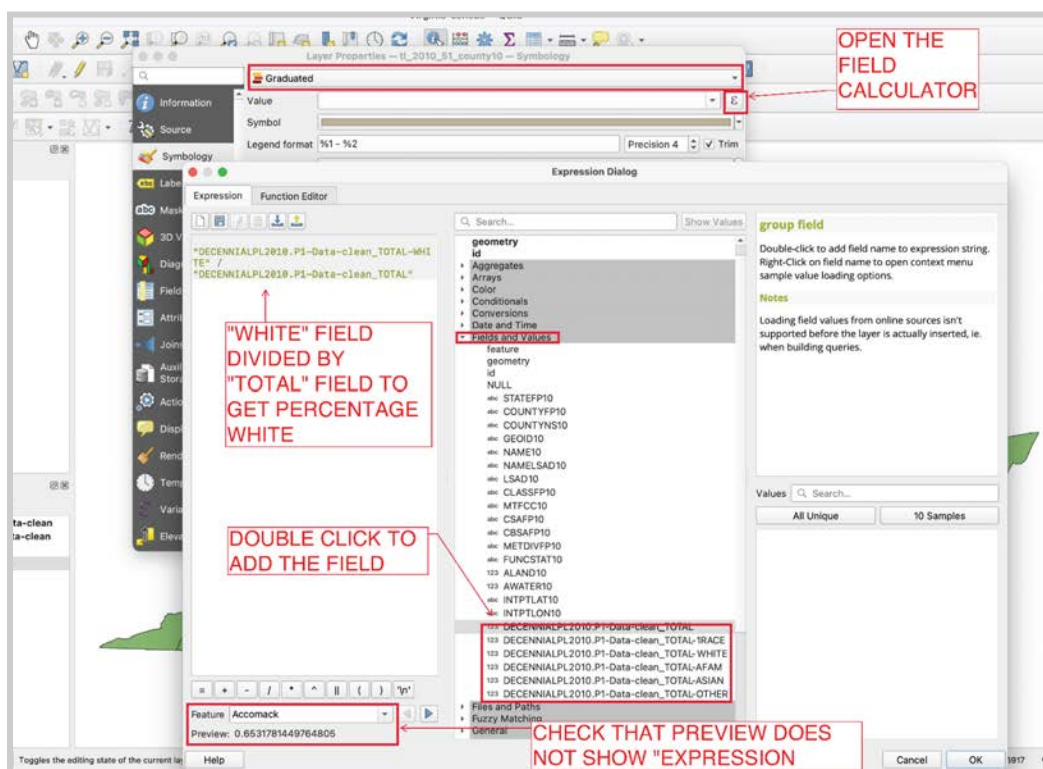
3d **Repeat** for the 2020 shapefile.

**Step 4: Color map by percentage of white and nonwhite residents.**

4a To get the percentage of white population, we need to divide the total number of white residents by the total number of residents (WHITE and TOTAL columns). We can do this using the field calculator. This field calculator allows you to manipulate one or multiple fields using algebra, but it also allows you to perform more sophisticated functions like converting text to numbers or using conditional statements (IF, THEN).
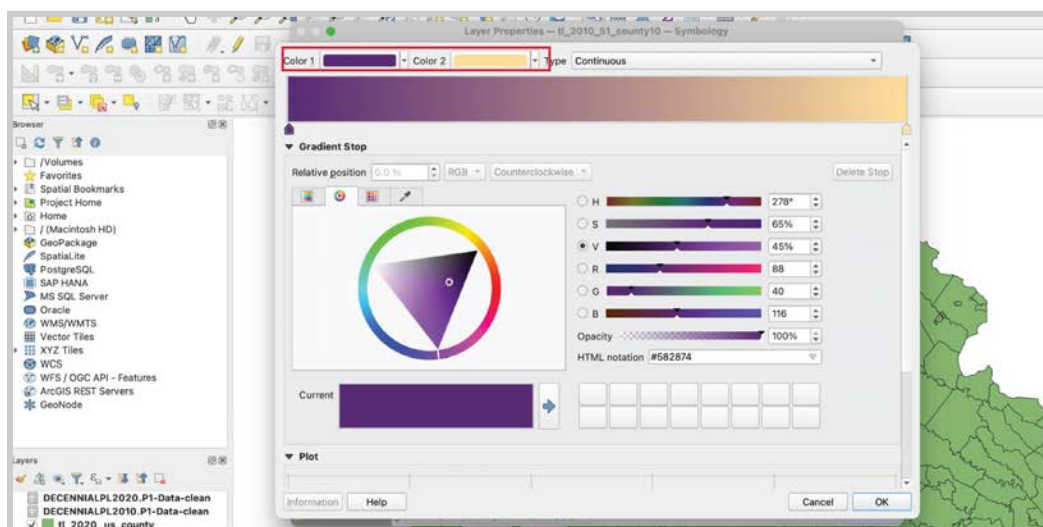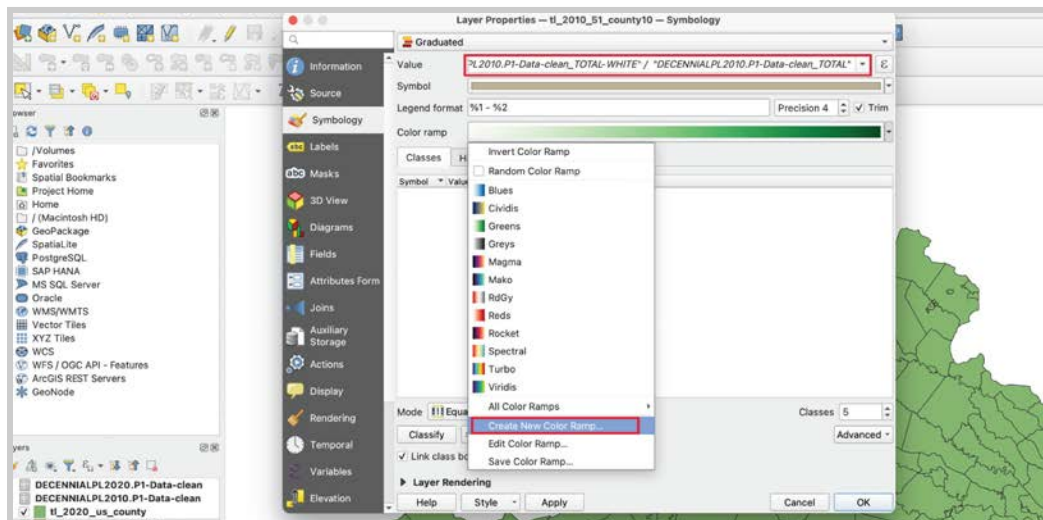
To open the field calculator, go to **Symbology** of your 2010 shapefile and select "**Graduated**". Click the **sigma (Σ symbol)** to the left of the "Value" field to bring up the Expression Dialog window. Here, you'll see all the possible functions in drop-down lists, as well as the "Fields and Values" of your layer. This is a useful way to easily see and add the fields in your Attribute Table without having to look up their exact wording every time. You can simply open up the "Fields and Values" section and double click to add your fields.

Add the "**White**" field, the division symbol (**/**), and then your "**Total**" field to get the percentage. You'll see it as a decimal number, which you can later convert manually to a percentage in your legend. To make sure you entered the formula correctly, check the "**Preview**" box near the bottom. If you see "Invalid Expression", re-check your formula.



NOTE: If you want to avoid the manual conversion, you can make the function produce a round percentage number by adding a couple elements: ROUND ( ("White-field-name" / "Total-field-name" )*100, 0), where 0 refers to the number of decimal places to round to.

4b Rather than a preset gradient, I prefer to make my own for this variable. I'm going to create a purple to orange gradient. In "Color ramp", select "Create new Color Ramp…". Then, you can set the Color 1 and Color 2 to your liking.

4c To show the data most accurately, you'll need to **manually** edit the values. Set your classes to "10", and then adjust the upper and lower values of each class bou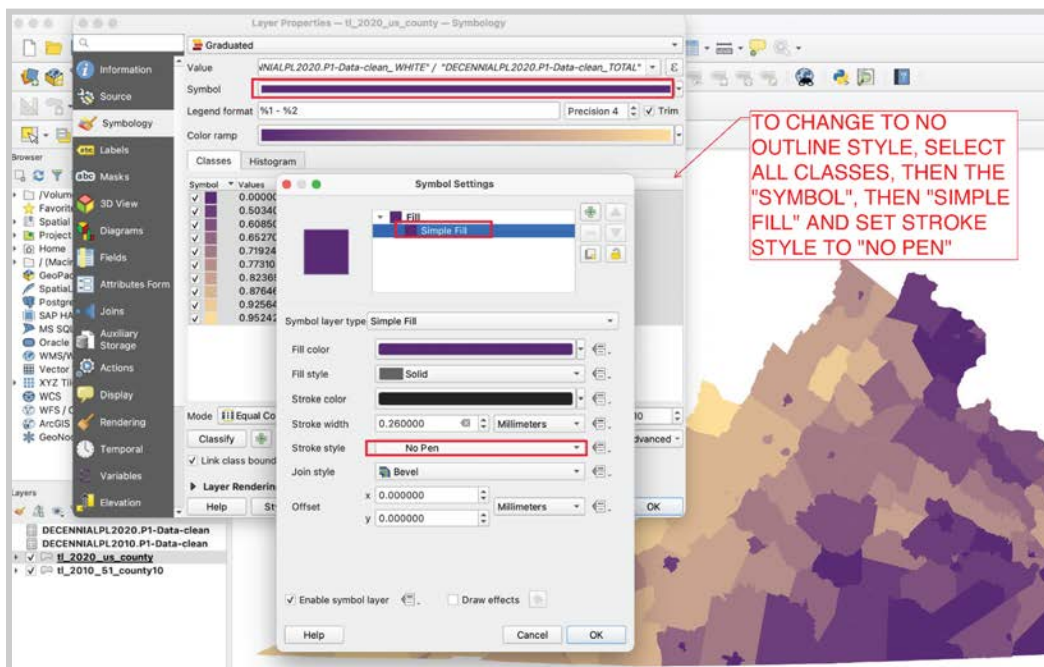nd to be a 10-percentage point range (eg 0-0.1, 0.1-0.2, 0.2-0.3, and so on – or, 0-10, 10-20, etc if you used the "Round" note above). To edit the class bounds, double click on the numbers in the "**Values**" column beside each Symbol.



4d If your symbols have an outline, consider removing them to see the data patterns more clearly, especially on smaller polygons. Select all of your classes and click on the "Symbol" dropdown near the top. Select "Simple Fill", and then dropdown "Stroke style" to "**No Pen**".
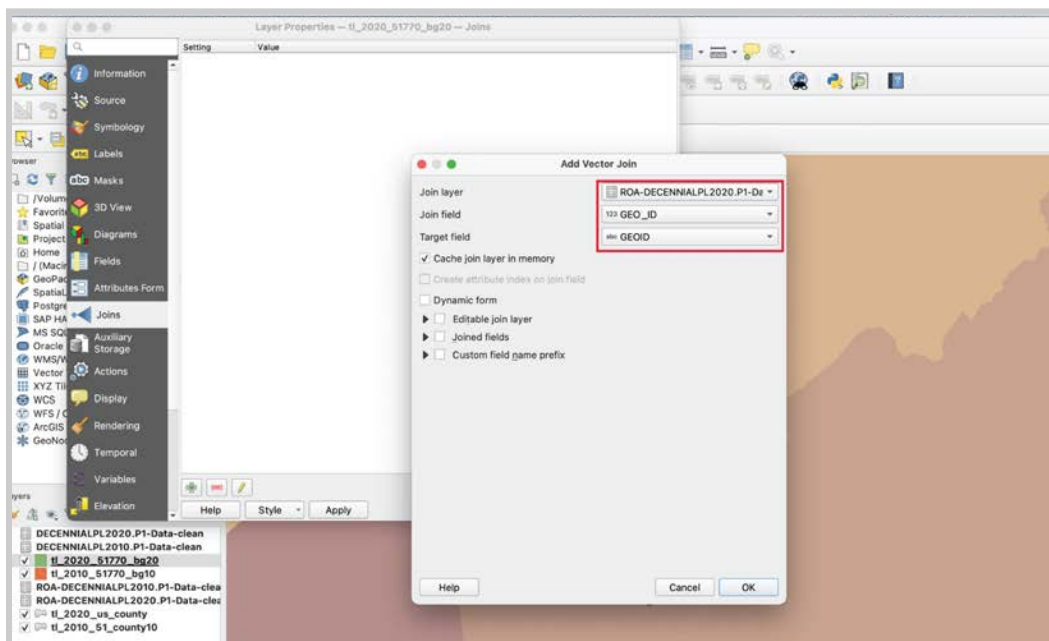


4e **Save** your layer style to use on the other layers, so that you can easily compare the data.

4f **Load** your saved style for the 2020 layer. NOTE that you will need to **update the fields in the field calculator,** since their names will be different ("2010" instead of "2020").

**Step 5: Add the city-level census data, join, and style.**

5a Using the previous steps as a reference, import the city-level block group shapefiles and csvs. Remember to choose the "DECENNIALP1" csv files for this tutorial. Join the 2010 to 2010 and 2020 to 2020, making sure to **join from the shapefile Layer Properties**, and not the csv layer properties.





5b Style these layers using the same saved style. Again, you will need to update the fields in the field calculator.

**Step 6: Add context physical geography layers.**

6a Similar to previous tutorials, we're going to add in a few layers to contextualize the data. I added in the Roads and River layers. You might also add in the Railroad layer, which will often draw the racial divide of a city.

Since we'll be showing the state and city levels both, we want to turn off the city layer when we show the state layer and vice-versa. We also don't need to show the roads and river when we're looking at the state level. To facilitate this, I grouped the shapefiles, data tables, and other geographies into group 1 (city-level census files, roads, and river), and group 2 (state-level information).
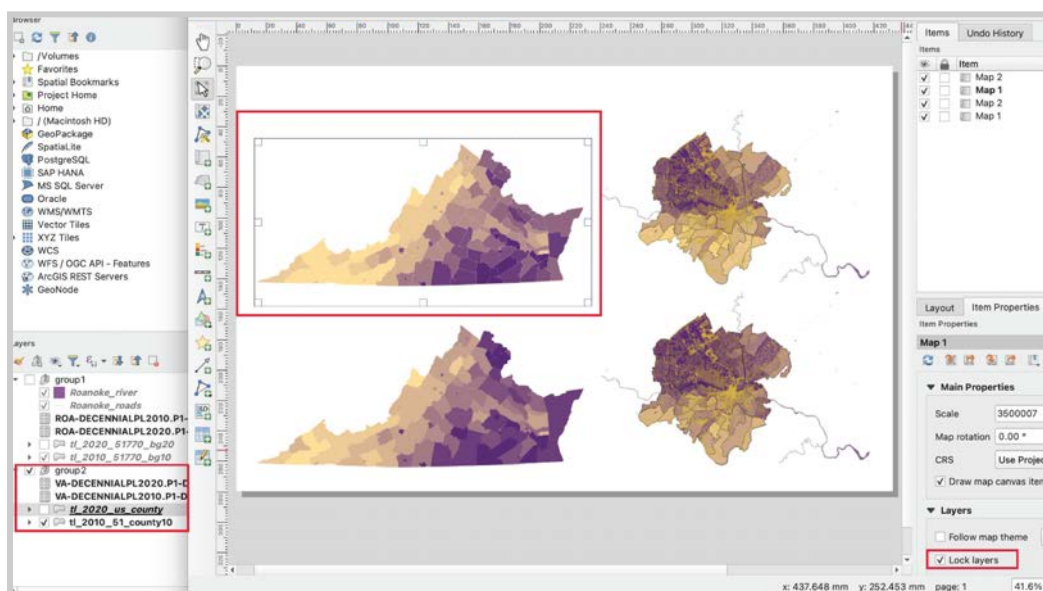
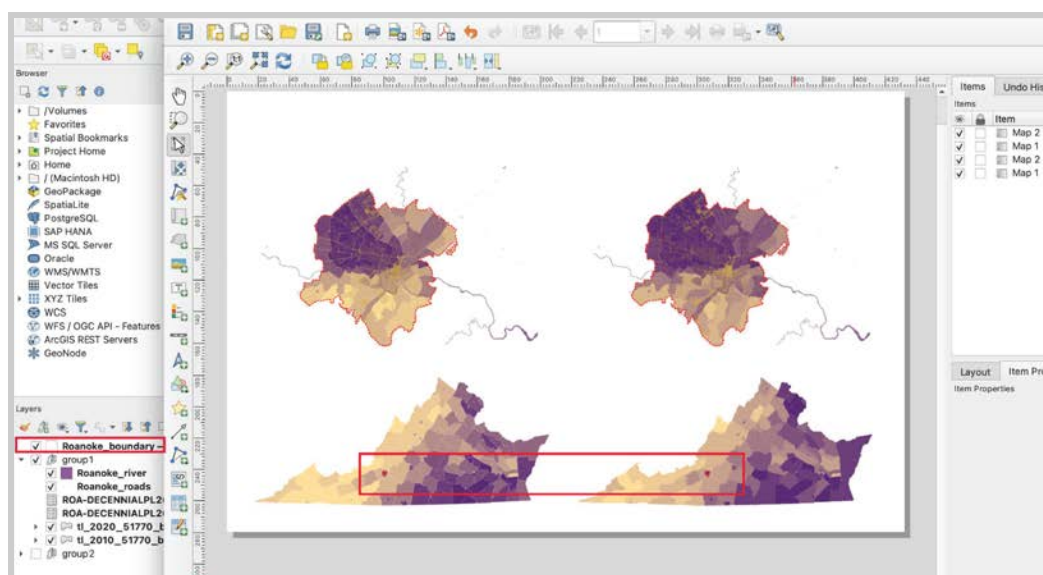**Step 7: Create print layout to compare the maps and years.**

7a Set your page to 11"x17".

7b Next, add a map of the state of Virginia. Turn off your city layers and your 2020 state-level layer. Check **"Lock Layers".** Copy the state map, turn on 2020, and then uncheck and re-check "lock layers" to update the second state map to show the 2020 data.
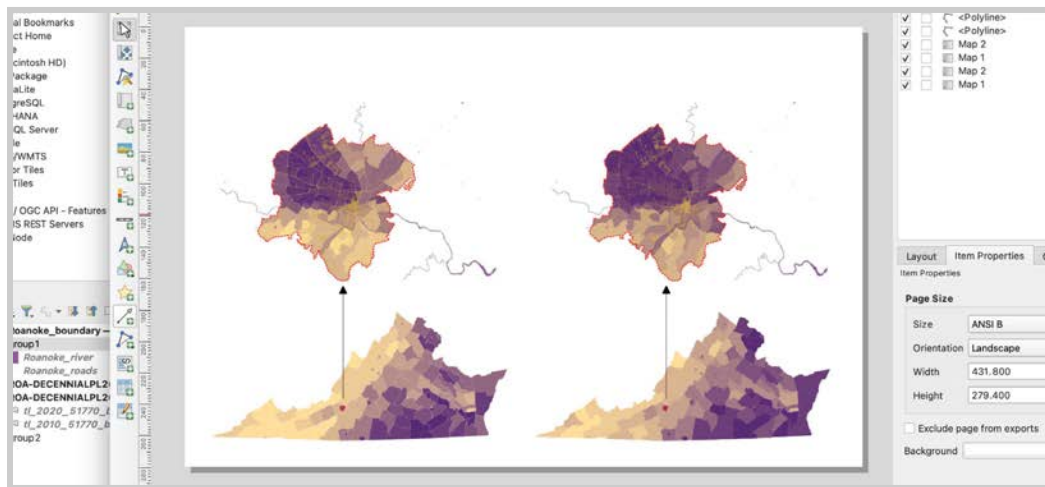
7c Do the same for the city-level data, turning off the state-level data and zooming in.
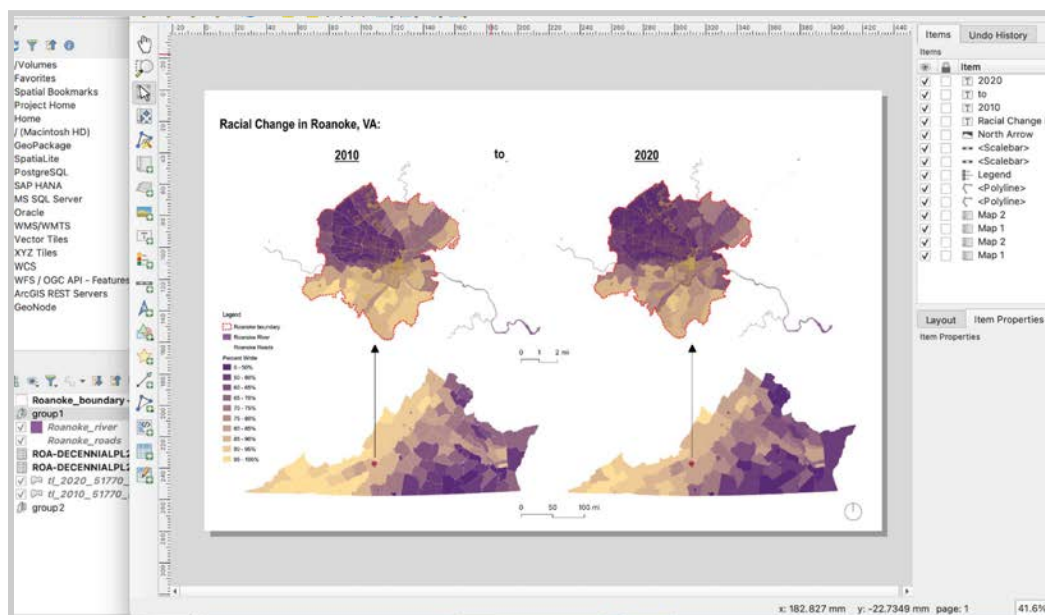


7d Now we can see the state-wide change in data between 2010 and 2020 as well as in the specific city, but without linking the two the information is not as legible. To make it clear what part of the state you're zooming in to, add the "boundary" layer to your map. Color it something bright like dashed red outline, and put it above the other layers so it always shows up.

Unlock and re-lock all your maps, turning off and on the relevant layers. You should now see the boundary clearly indicating the city's location on the Virginia map. Add an arrow from the city in the state map to the blown-up city map, to make things even clearer.



7d Arrange your maps by year. Add a legend (be sure to indicate units, in the case "percentage"), north arrow, and title / subtitle. Add a byline and sources. Add separate scales for the city-level and state-level maps.
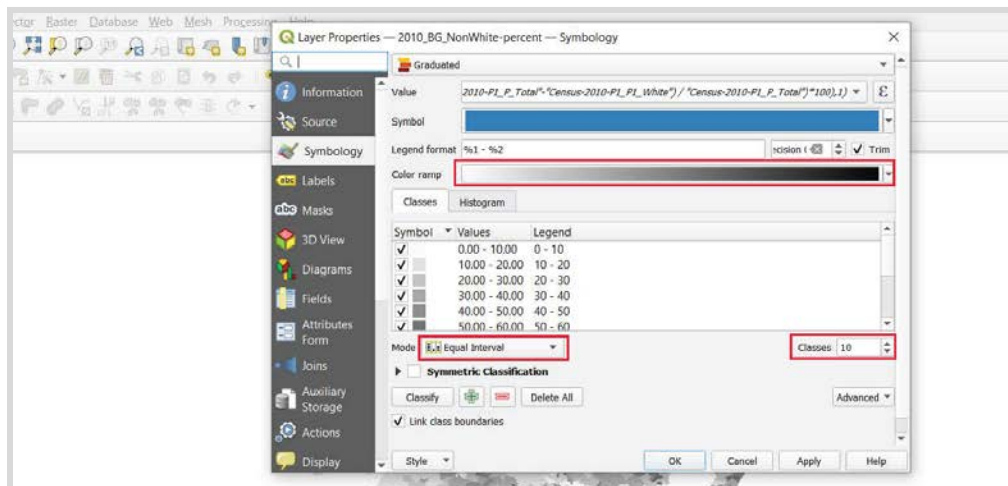
- **Bonus (STRONGLY RECOMMENDED) -**

**Step 8: Raster analysis.**

In the previous steps, you visually analyzed the difference between 2010 and 2020. However, it's difficult to see the difference without flipping between the layers. You could create a GIF to make the change more evident, but you can also perform raster analysis on the two images. QGIS has a tool called the Raster Calculator that allows you to perform math on the pixels in images. In other words, you can export the two years of census data and use simple subtraction to visualize the change between them as a new image.
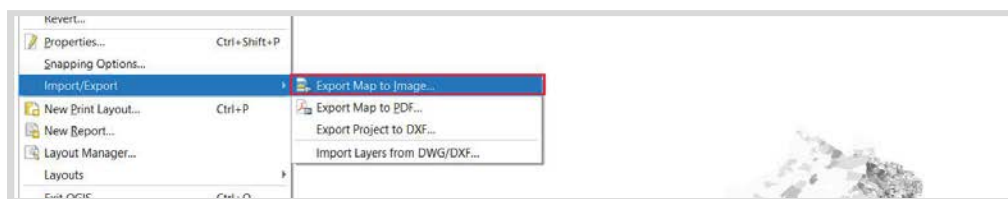
This is especially useful when comparing the changes between census decades, since their shapefile outlines have changed. Because of this, you can't get an entirely accurate picture just by comparing their spreadsheet tables. The Raster Calculator will compare them pixel by pixel, rather than using attribute tables to compare them block ID by block ID.

> **NOTE:** to accurately compare the two datasets, first make sure that they're styled identically. In other words, both maps should have the same color scheme and the same percentage ranges (use "Equal Interval" in the graduated styling mode and set the same number of "Classes").
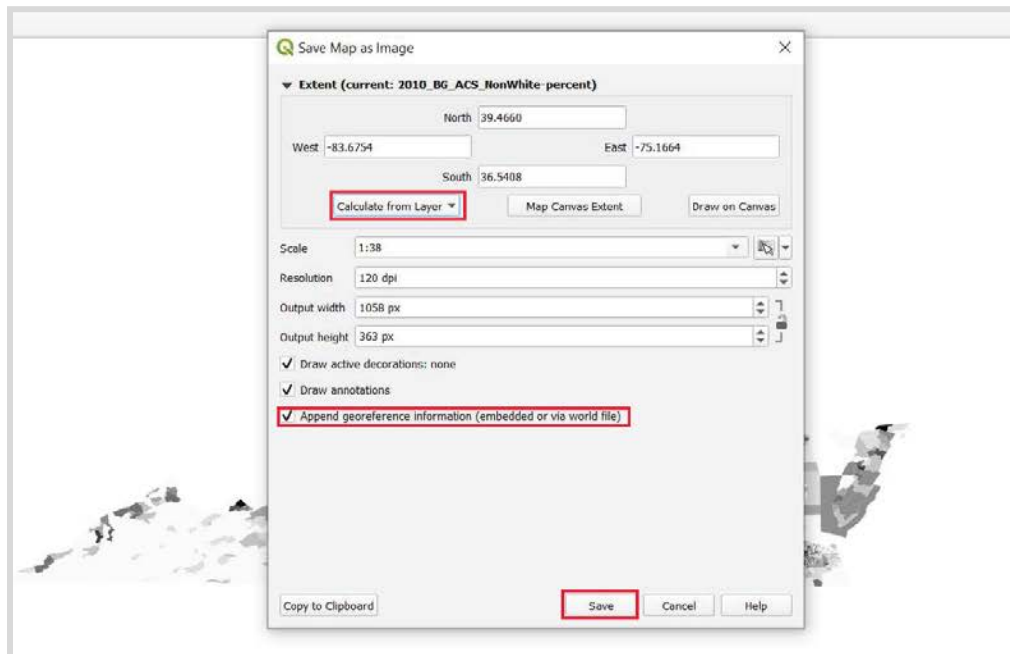
8a First, style your 2010 and 2020 maps with a **black and white gradient** rather than a color gradient. This will make the Raster Calculation easier. Double check that your percentage ranges match between the two layers.
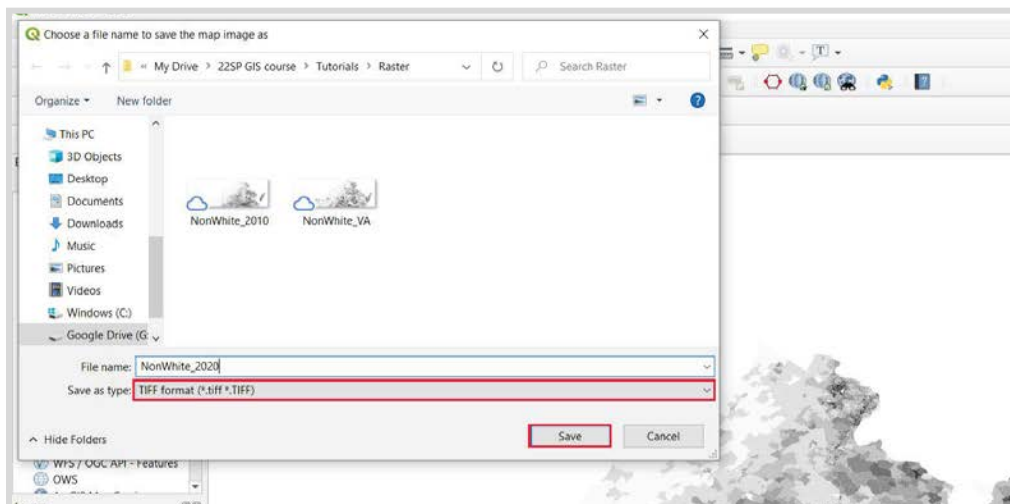


8b Next, export both the 2010 and 2020 maps as images. Turn off all your other layers, and only leave on the 2010 map. Go to **Project > Import/Export > Export Map to Image…**

In the dialogue box, set the **Extent with "Calculate by Layer",** and select your 2010 layer. Leave the other fields as their defaults; make sure that "Append georeference information" is checked.
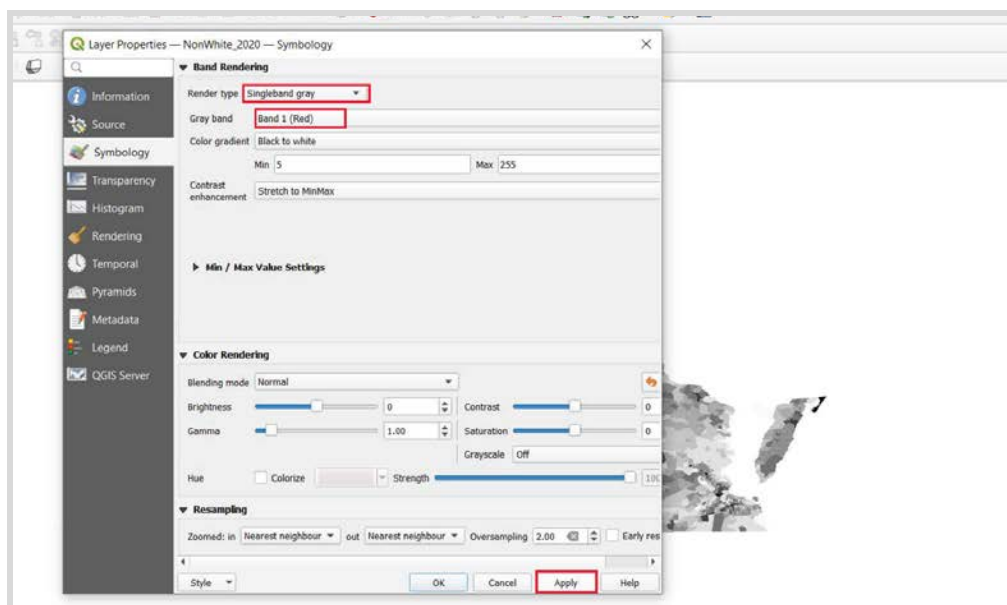


Click Save, select **"TIFF",** and save your new Raster image somewhere safe. Repeat this step for the 2020 layer.



8c Add in your new Raster layers with Layer > Add Layer > Add Raster Layer…

8d Double click on one of your census raster layers. Set the property "Band Rendering" (the first property under "Symbology") to "**Singleband Grey**". You can leave the "Grey Band" to "Band 1 (Red)". Make a note of this, since we'll want to compare the Band 1 of both layers in the next step. Select "Apply"
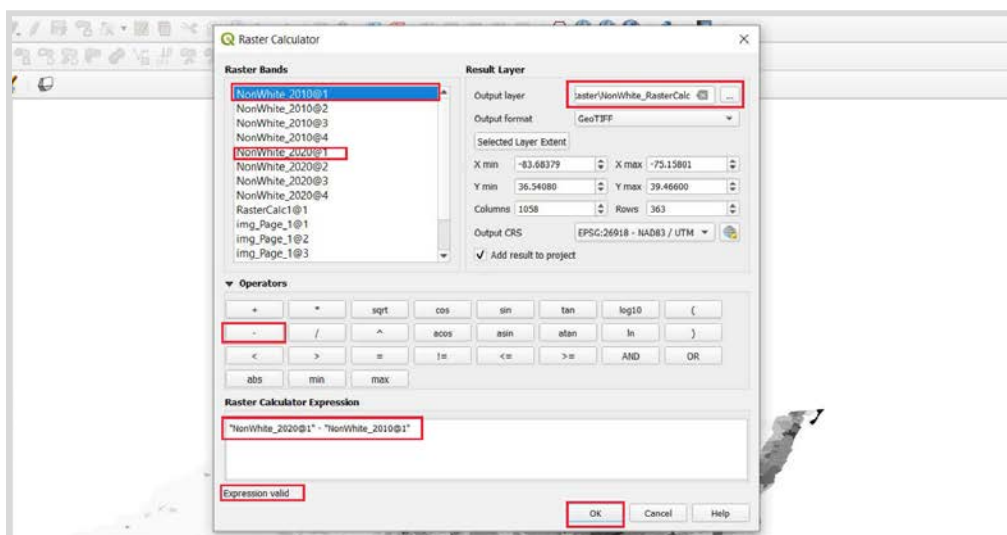
Repeat this for both layers.

8e Now, open **Raster > Raster Calculator**. Since we can see that the Nonwhite percentage has grown between 2010 and 2020, we want to subtract 2010 from 2020. As with the Graduated Symbology expression editor, you can double click to add layers to your Raster Calculator expression. First, Double click on the 2020 Band 1 layer listed in the "Raster Bands" list (upper left). Then click the minus (-) symbol in "Operators".
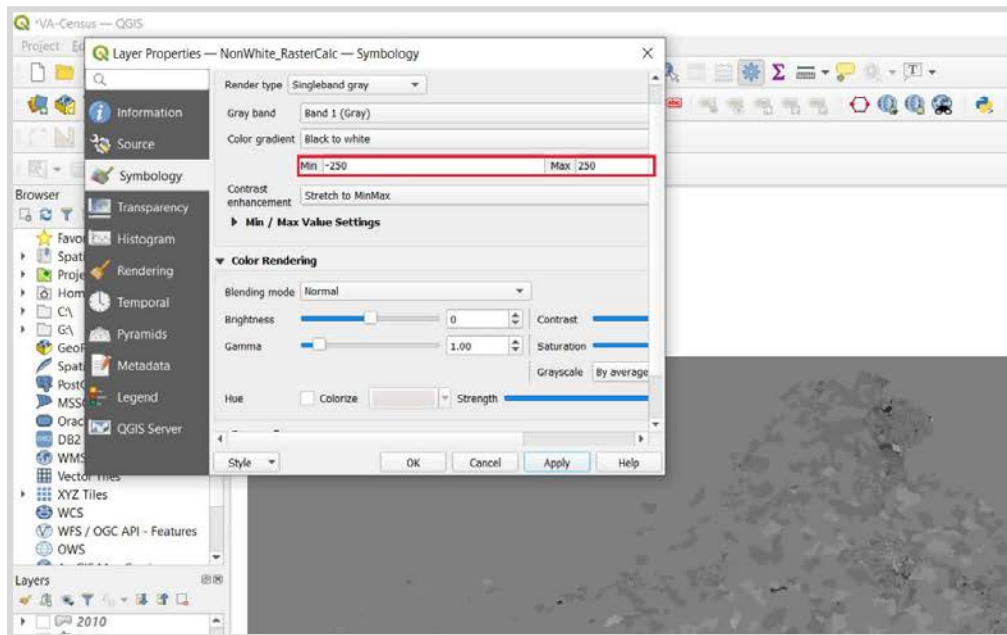
> **Note:** Notice that at this point, before completing the expression, you have a text message in the bottom left that says **"Expression invalid"**. This text, which appears in all QGIS expression editors, is a useful indicator to keep an eye on.

Now, double click on your 2010 Band 1 layer in the "Raster Bands" list to complete the expression. In the upper right, set the "Output Layer" to a safe place with a logical name. Click "OK".
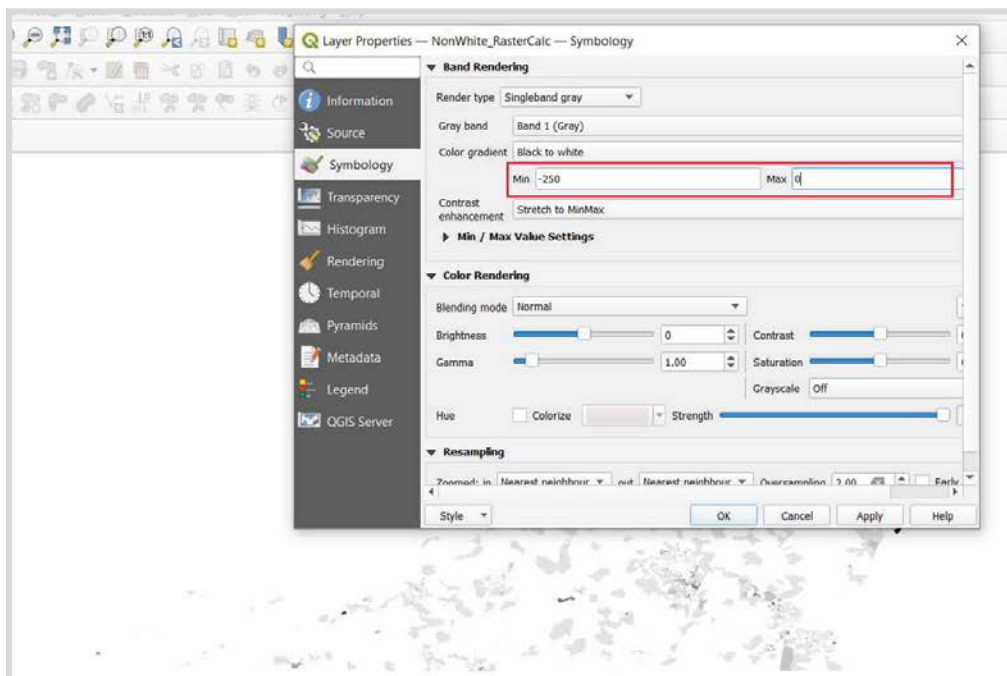


8f You should now see the new Raster Calculated layer added to your map. It will automatically set its color scale from -250 to 250. The grays that are LIGHTER than the background gray are negative values (eg areas where

NonWhite percentage DECREASED from 2010 to 2020); the grays that are DARKER than the background are positive values (where NonWhite percentage INCREASED).



If you only want to see the decrease, double click on the layer and set your Band Rendering MIN to 0. Leave your MAX at 250. Likewise, if you only want to see the increase, set the MIN to -250 and the MAX to 0.



8g Add a page to your print layout with this new map. Title it something like "change in nonwhite population 2010 to 2020"