

A User-Friendly U.S. Census Browser for R

Kiegan Rice

1 Introduction

Census data is an important snapshot of information about a country at different times throughout their history. It is an integral part of keeping track of the history of a country and the people who occupy it, as it gives us records of the occupants of a country and how they live their lives. The value of census data to a country and its citizens is difficult to overstate. Many recognize the usefulness of the census data, especially when aggregated and presented in a way that allows for viewers to learn something new about the world around them. While the history books we learn from in high school present us with a narrative of the events that occurred, students often don't get to interact with the raw data ourselves in that learning environment. Accessible and usable census data allows the exploration of different demographic groups over time, or investigations of a particular period of time and what the demographic and economic landscape looked like in the past. Even today, for those who can't travel the world, data about the world around them could allow them to learn more about those places and learn something more about groups of people they may not usually engage with.

From an early point in the United States' history, there were "many eminent men of science" who recognized the value of the census data and worked to aggregate and present the data that had been collected on the population. Francis A. Walker's "Statistical Atlas of the United States", based on the 1870 census, was an impressive effort in aggregating population data to present it in a visually appealing way. Although the Census Bureau's "Statistical Atlases" eventually stopped being made, they were an important start to the effort of presenting census data to a wider public.

Today, as methods of data analysis and visualization continue to be developed and improved, access to census data allows us to look back on that history and explore, synthesize, and visualize the information to learn more about patterns in many different parts of the population.

Mention Dr. Hofmann's paper to improve the Statistical Atlas here... forthcoming? published? As well as Haley's ggmosaic stuff

Ever-improving visualization and data-wrangling methods in R give those interested in statistical graphics a wealth of abilities to explore and learn from data; however, it is difficult to make use of these tools on census data if that data is not available and easy to explore in one location.

An inherent problem in census data is that a country's census changes over time; the variables collected, how they are collected, and even the locations they are collected for change as the country is formed, and subsequently grows and changes. The United States census data is no exception to this rule. In a little under two and a half centuries, the census has taken on many different forms. Data on occupations has transformed as the employment landscape has changed; new states have been formed, the most recent being within the last 100 years; definitions of various demographic groups and the terminology used to describe them have been updated as the demographic makeup of the country has changed. Each decennial census brings a different set of variables to the table. Sometimes these variables are new things the Census Bureau is interested in learning, while sometimes they remove variables that are no longer relevant or whose information is captured somewhere else.

Unfortunately, because the founders of the U.S. Census were unable see 200 years into the future, those interested in working with census data are left with quite the inescapable mess. If you want to focus in on a particular demographic group and their journey as part of the population of the United States, you may have ten or more different variable names to describe that one group over the course of the census from 1790 to 1960 - and that is just for one single group! Of course, we cannot just simply change variable names to match our own research needs. It is important to keep the data in its true form and be honest to the way that the population was defined at different times throughout history, even if our instinct may be to ‘clean’ the data by changing variable names.

This, of course, leaves the user with a wide variety of variables that are far from consistent across years. In order to track one demographic group across years - let alone many groups - a clean user interface that helps users see exactly what information is available to them across years is a necessity. To streamline the process and assist researchers in finding out what information they have access to and what they don’t, we present a U.S. Census Browser for R, with the user interface being a Shiny application, and the downloadable files being ‘tidy’ csv files that those with a small amount of R experience should be able to work with.

2 Background ?(Literature Review)?

The University of Virginia Library hosted a “Historical Census Browser” for many years that allowed users to search United States Decennial Census Data for use in research, teaching, and personal inquiry. The data included data on various aspects of the U.S. population from the 1790 Census through the 1960 Census, originally populated using the ICPSR 3 dataset. This Historical Census Browser was free and available for use for anyone with an internet connection. This Census Browser allowed a user to peruse available topics for each census year, at both the state and county levels. The Historical Census Browser was taken down on December 31, 2016, with the county-level aggregated data becoming unavailable several months before this. The source was widely used, with several other university libraries and educational resources including University of California Santa Barbara, University of Pennsylvania, University of Michigan(University of Michigan Population Studies Center 2017), and the Smithsonian’s History Explorer directing researchers to use the University of Virginia source.

The University of Virginia Libraries website now directs users to “Social Explorer” or the “National Historical Geographic Information System” (NHGIS) website. Social Explorer, which is populated with the ICPSR 2896 data, requires that users pay to use (or use through library access from a library who has paid for it), and does not offer full download. NHGIS, hosted through the University of Minnesota, is very difficult for users to navigate when looking for specific information across multiple years. The Institute for Social Research at the University of Michigan (ICPSR), has the full ICPSR 2896 data set, split into 106 separate data sets, but requires that users be a part of a member institution, and requires users to agree that they will not distribute the information in any way after gaining access to it. It also does not have any browser function for users to look at the data unless they download the entire file for each year, in either a SAS, SPSS, ASCII, or Stata set-up.

None of the aforementioned resources provide the full advantages that the Historical Census Browser offered: a free-to-use and user-friendly data browser that allows users to choose which data they are interested in using, and download the complete (aggregated) records for their own independent use. Although the Historical Census Browser data is now dated based on the existence of an updated,

cross-checked version of the decennial census through ICPSR (2896), it is for the most part accurate and allows for free and open use for researchers and others alike.

3 Work

3.1 Data

Although the county-level data had already been removed from the website, the state-level aggregated records for each decennial census from 1790 through 1960 were captured from the website in September of 2016 and saved as raw data with the intent to create a resource for R users that allowed the same main functionalities of the original University of Virginia Historical Census Browser, streamlined for easy data searches and data management.

The majority of files for the decennial censuses (sp?) were complete upon capture from the website. However, there were two years of census data that are missing some column names and thus are unverified data. For the 1890 and 1920 censuses (sp?), columns were compared to all columns from ICPSR 02896 5 data files by both correlation and Euclidean distance, and thus some columns were able to be correctly identified.

Describe this process in more detail once you actually get this shit figured out. How many columns were correctly matched? How many weren't? Were we able to at least get a foothold on which columns go where in the excel spreadsheet after identifying a couple of them?

3.2 Description

- What the package includes and what capabilities the Shiny app has

3.3 Example

In the Shiny app, we can choose to focus in on a single year of the census or look for data across a range of years. We can easily use the “Get Your Data” Shiny app to search for a topic of interest, find the data we want, download state-level information, and tell a visual story of populations in the United States over time. To demonstrate the all-knowing power of this Shiny app, we will walk through an example on the history of the African American population in the United States. We begin by running the Shiny app:

```
library(shiny)
shiny::runGitHub("kiegan/censusbrowseR", subdir = "shiny")
```

The default set of years is 1790 to 1920. However, we can easily expand the slider range to be able to explore data across all the available years, 1790 to 1960. We can also order by how the “Total” column, or how many of years each of the variables appear in. These variables appear on the left side of the table, as rows.

The available variable names here are already ordered by how many years they appear in, so we can see right off the bat that we have a serious lack of variable continuity. The identifying variables that appear each year are YEAR, STATE, TOTAL.POPULATION, and TYPE. That means the total number of times they appear is 10, the number of census years we are looking at. After that, the next most

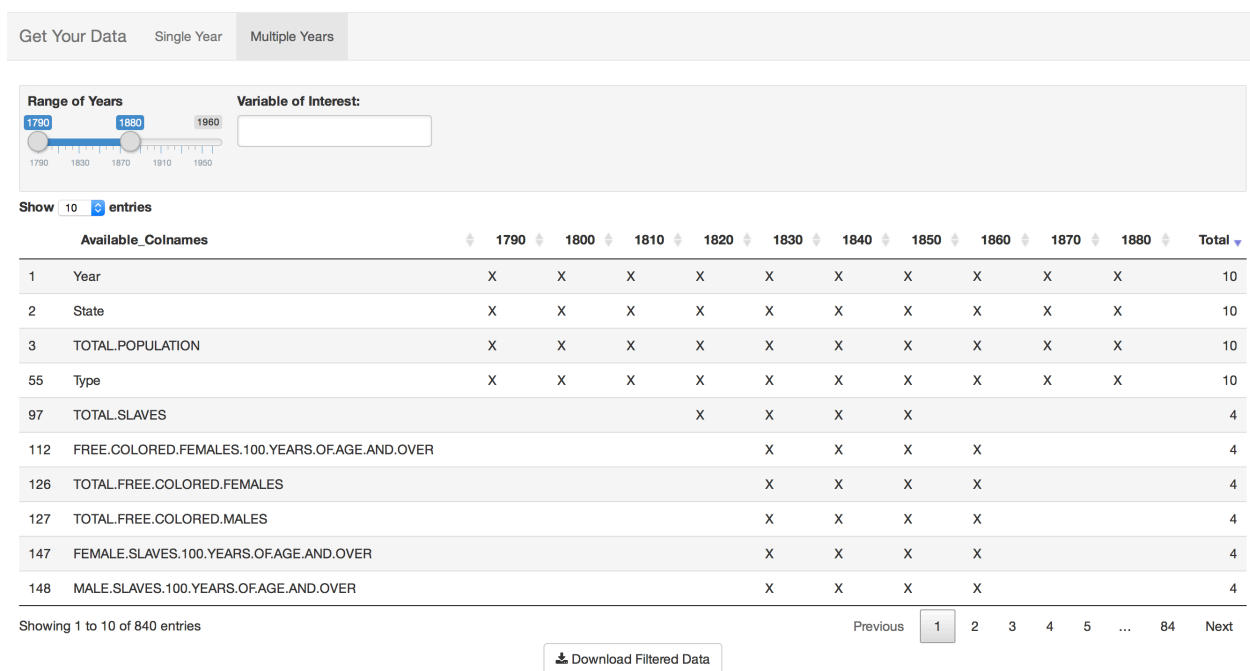


Figure 1: An example of available column names for the years 1790 to 1880.

common variable only appears in 4 of the 10 years we have chosen. That means we are up against some pretty significant hurdles to track any demographic group over time.

For now, we are focusing in on African Americans throughout U.S. history, which means we need to start with the term **SLAVE** in the early years of the U.S. census. The vast majority of African Americans were slaves when the United States was founded, so variables that count the number of slaves are important variables to have in order to get a grasp on the African American population in early U.S. history.

Once we have our range of years chosen, we simply enter a search term. Searching for **SLAVE** gives us a list of all the possible variable names across all years, and we can again sort by how many years each variable name appears in. Users are then able to click each of the variables they are interested in, and select **Download Data** in order to download a .csv file with the data for all states in the selected years. This .csv file will also always include the columns **YEAR**, **STATE**, **TOTAL.POPULATION**, and **TYPE**, even if they are not selected by the user.

For this example, I executed three separate searches within the **Get Your Data** app, using three different terms that the census used to categorize African Americans at different points in the decennial census. This gave me three separate .csv files, which I was able to easily combine using `full_join` from the `dplyr` package. Now, for an entire demographic group, all of the state-aggregated population counts are together in one data frame, although different terms are used at different points in time.

We can now begin to explore this data. We can very easily plot state-level counts of the **SLAVES** variable for the year 1790. We recommend the use of the **USAboundaries** package in concert with this census browser, as it provides the most accurate boundaries of the United States at any given date, and it is important to realize that values in the data are for the states *as they were* during that census, not the current (2017) boundaries of the states.

We now have a visualization of the slave population in the United States in 1790, and we can easily look at subsequent years of the decennial census to see how the population changes over time.

If we jump to the year 1820, the census now includes a count of **TOTAL.FREE.COLORED.PERSONS**, as now there are citizens who are African American or other minorities that are free, and not slaves. (yay!!!!)

Moving forward from there, we can start to see interesting patterns once slavery is abolished. In the year 1850, we can also start investigating what the total African American population looks like in each state, rather than just the slave population and free population.

3.4 Implementation

How to use it - this is the most technical part

4 Discussion of Future Work

- County-level data would be great
- Current data (up to 2010 Census)

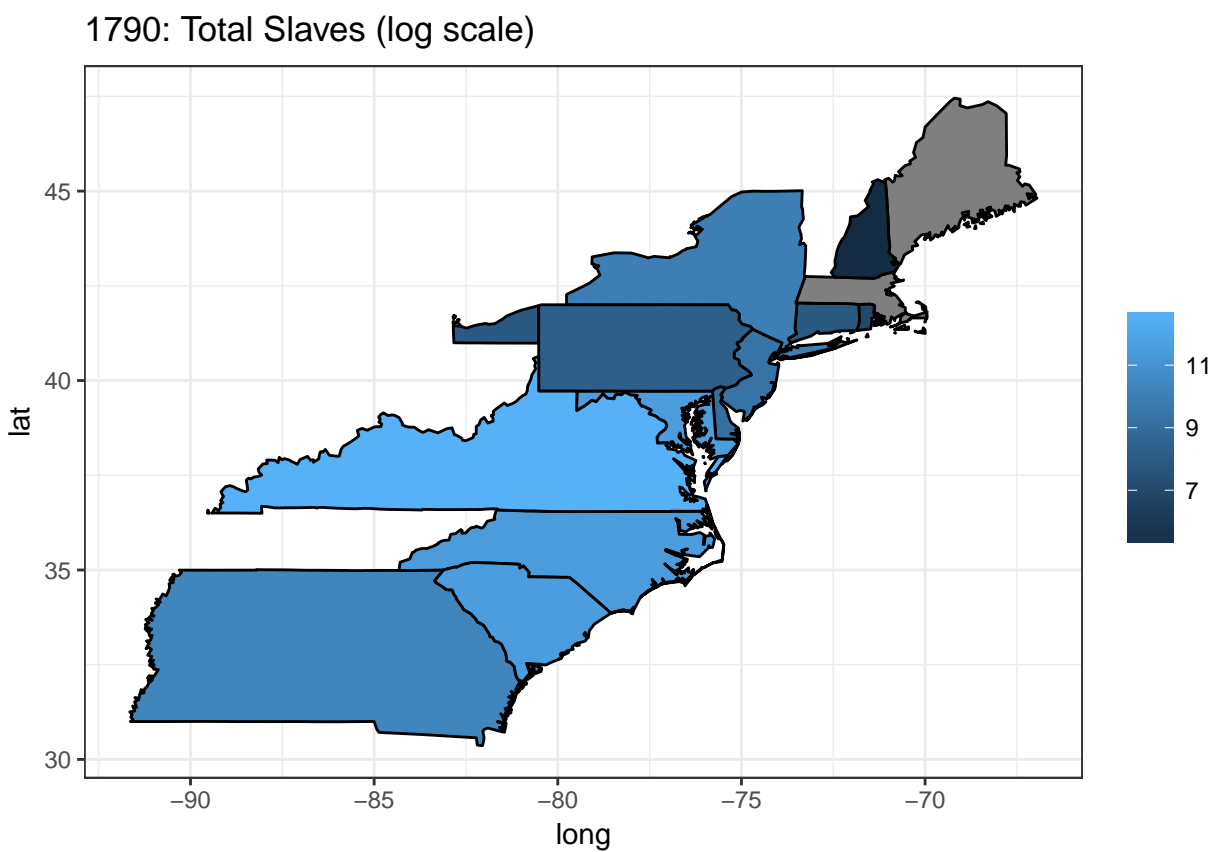


Figure 2: Total number of slaves per state in 1790, plotted on a continuous log scale. State boundaries for July 4, 1790 were gathered from `USAboundaries` package.

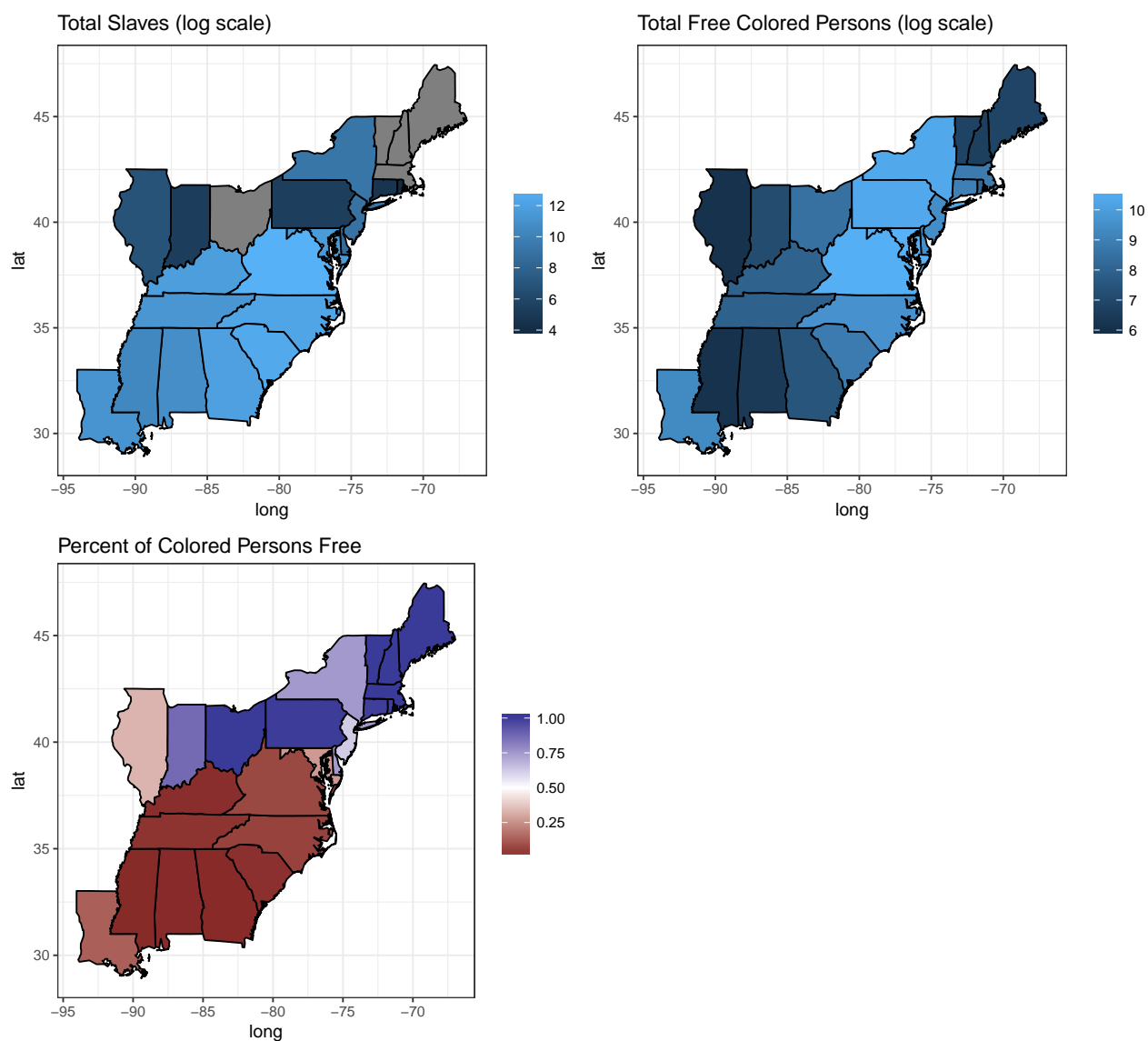


Figure 3: Number of Slaves, Number of Free Colored Persons, and Percentage of Free Colored Persons in 1820.

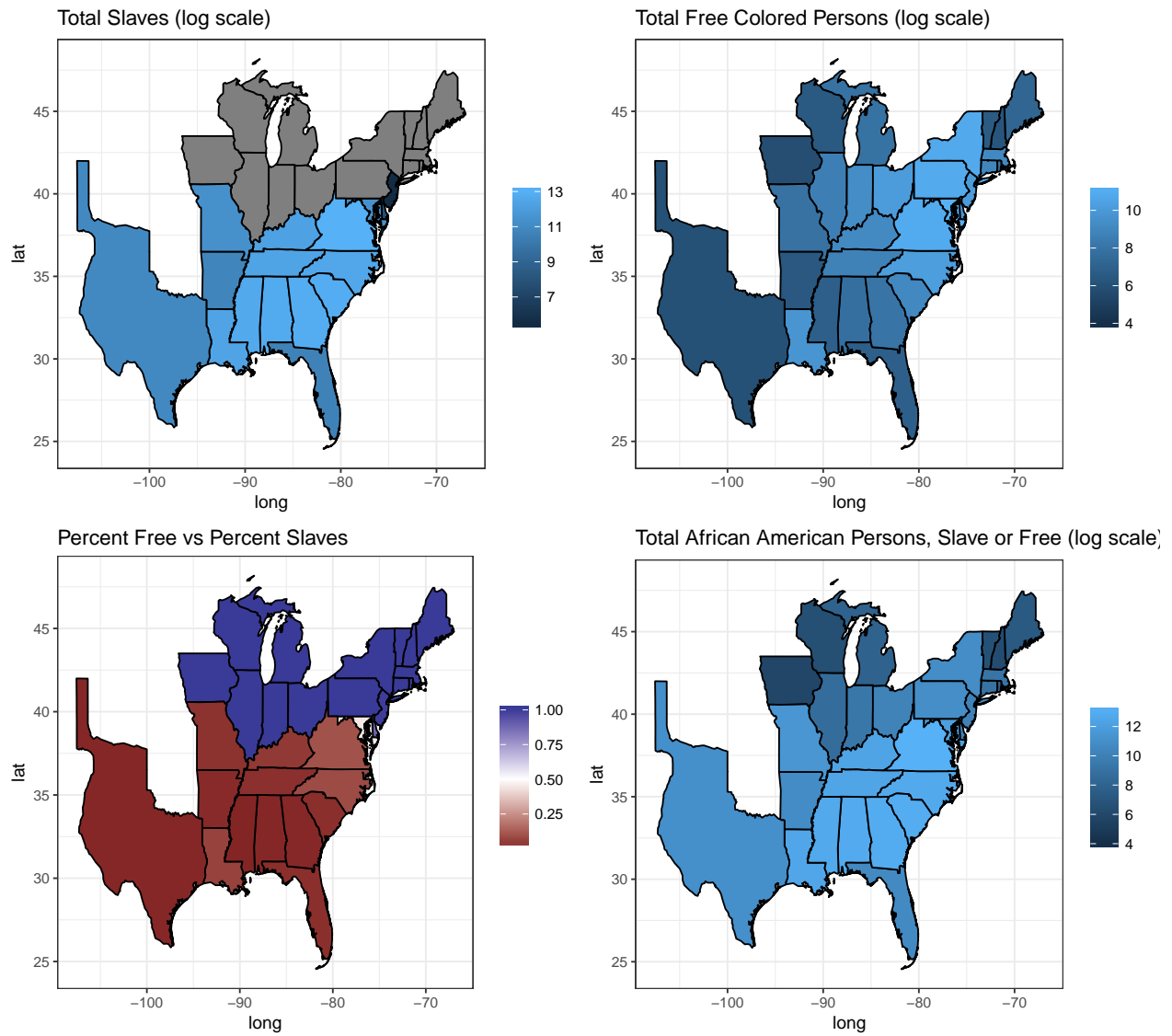


Figure 4: Number of Slaves, Number of Free Colored Persons, Percentage of Free Colored Persons, and Total African American Persons in 1850.

5 References

University of Michigan Population Studies Center. 2017. “Historical Census Browser (1790 - 1960).” Accessed March 27. <http://www.psc.isr.umich.edu/dis/data/resource/detail/1369>.