

# Revisiting Du Bois

## A Data Visualization Project

Nathan Kim

Advised by Professor Elisa Celis

## Introduction

Prolific and prominent Black writer, historian, professor, political activist, W.E.B. Du Bois has left a legend in many ways. In addition to being a sociologist who meticulously drew and recorded social statistics for over fifty years, Du Bois was also a Black theorist of knowledge, an educator, a novel writer, and a historian. Some of these perspectives are visible in his approach to data visualization. While other data visualization pioneers like Edward Tufte try to objectively gauge the value of a dataset as the proportion to which it represents “truth,” Du Bois recognizes visualization as a creative endeavor alongside photography and historical analysis.

My project seeks to explore his work with a critical and quantitative lens. I hope to bring to Du Bois’ portraits the worlds of modern statistics, computing resources, and clearly defined systems of visualization, but also to these fields bring Du Bois’ humanist, expressive lens. To do so, I began development on the `ggdubois` graphics library for data visualization in the R programming language. In extending Du Bois’ data visualization through practical tools, I hope to also convey his artistic perspective on data visualization.

# Background

## A short history of data visualization

Modern data visualization began with modernity itself. The map as a medium was developed first, alongside modern notions of states and borders throughout the fifteenth and sixteenth centuries.<sup>1</sup> States had new reasons to formally define borders through a codified text and recognizable medium – the map. Other forms of data visualization, like bar and line charts, were developed over time to similarly convey graphically what language had to stretch to express. The earliest known use of these graphics to convey statistical information appeared in 1644, when a Flemish astronomer named Michael Florent was tasked with representing several measurements of the distance between Florence and Rome. A table might have sufficed to convey these distances, but the graphic can also represent the variation and the relationship of these distances to each other and easily show any outliers.<sup>2</sup>

The next three hundred years brought much change to what is known as data visualization, but the idea of translating statistical concepts in graphical ideas remained the same. Most fundamentally, the idea of *measurement* to produce a graphic solidified during this time, assisted by both new tools and new schools of thought to guide them. As the age of Exploration, colonization, and the Enlightenment as a whole took form, so too did units of distance and time to measure and quantify colonizers' discoveries. First there emerged mediums to quantify raw distances and quantities, like maps for navigation, but soon followed more elaborate and abstracted forms of describing space, time, and everything in between. Topographical maps involving contour plots emerged in 1584 but became more common in the 1700s. The first known scatter plot emerged in 1833. John Snow in 1854 created his famous maps of the London cholera outbreaks, at a time when the map as a medium was familiar but the application of such a medium to convey statistical trends was not. The transformation of scales to accommodate non-uniform data became prominent in (1863, 1869). All of these methodological innovations were accompanied by technologies like color printing and logistical projects like nationwide censuses that cemented visualization as a canonical form of understanding data. Besides a brief period that some scholars call the “dark

---

<sup>1</sup>

<sup>2</sup>Friendly, “Milestones in the History of Data Visualization.”

ages” of modern data visualization, discoveries like these progressed steadily from 1600 onward.

Today, the accepted canon of data visualization pioneers has coalesced around Edward Tufte, Jacques Bertin, and John W. Tukey.<sup>3</sup> Bertin, in his *Semiology of Graphics*, established a conceptual backbone that related visual elements directly to trends within data through graphical components, an idea from which other graphical “grammars” (including ggplot2, as mentioned below) would one day develop. John Tukey, in his *The Future of Data Analysis* (1962) and *Exploratory Data Analysis* (1977), argued for data analysis as a branch of statistics distinct from mathematical statistics and argued that recognizable and reproducible visualizations were key parts of understanding data. Edward Tufte, in *The Visual Display of Quantitative Information* and many other texts, argued for the judgement of the quality of data graphics as how “truthfully” items in the graphic corresponded to measures in a dataset.

### ggplot2: a grammar of graphics

The perspectives of Tukey, Tufte, and Bertin are embodied today in software packages for data visualization. Data visualization is no longer a painstaking task of an artisan who must draw every point by hand; just as Tukey once argued for repeatable and easy-to-produce exploratory graphics, the simplest data graphics now can be drawn in a single line of code. Bertin and Tufte’s emphases on deconstructing graphics are embodied today in the ggplot2 library in the R programming language.

R has become popular today for many reasons. It is designed to be more user-friendly than other lower-level languages, automatically handling memory allocation and variable typing for the ease of the user.<sup>4</sup> Being designed specifically for statistical analysis, it has many native data structures and functions for computing models. For example, almost all data types in R are natively interpreted as some form of a vector, and functions for computing linear algebra are highly optimized in C++ before being ported to R. Likely the largest contributor for R’s popularity has been its open-sourced nature. As well as being free, a considerable advantage compared to languages like Stata and SAS that can cost hundreds or thousands

---

<sup>3</sup>There are of course many, many scholars that laid the foundation for modern data visualization besides these. Some examples are....

<sup>4</sup>The downside of this high-level design is that R’s performance has come under constant criticism. Countless Medium articles have been written on R’s performance against languages like Python, Julia, and C. This subject is left to a footnote as it is not the main subject of my paper, and the most important takeaway I feel is that despite any potential performance issues R is still a fairly popular language among data visualization professionals and statisticians.

of dollars a year, being open-sourced means that user extensions are at the core of R’s functionality.

One especially important extension in R is the `ggplot2` library and the ecosystem of user-contributed software packages it has spawned. The `ggplot2` package is popular for providing many utility functions for graphics in R, for example a function called `geom_smooth` for smoothed conditional means. These functions can often abstract away complex logic from users, so that (in the case of `geom_smooth`) a potentially complicated choice between loess and general aggression models can be hidden from the user that simply sees a best fit line.

But the library is far more influential for contributing its eponymous *grammar of graphics*, or an extendable logic for how to make plots. Plots in every programming language can quickly become complex and syntactically verbose because of how many geometries and aesthetics of a plot there are to consider. `ggplot2` was created as a response to this situation in R, creating a grammar from which almost any graphic could be created. In this grammar, every plot can be decomposed into a few ingredients:

1. A dataset
2. One or more layers of geometric objects, or “geoms”, along with mappings between data and features in these geometric objects.
3. A coordinate system to translate conceptual *x/y* or *radian/theta* values to pixel values.
4. A theme, or aesthetic styling that are not directly related to the data itself.

A demonstration of this logic is shown in Block 1. The plot is initialize with `ggplot()`, then an arbitrary number of geometries can be added with the `geom_*` syntax, a scale is used to adjust how items in the data are mapped to aesthetics, and finally a theme is added for styling. The strength of `ggplot2` lies in how the individual components of a graphic are separated by a plus sign (+), or in other words that visual components can be written directly through syntactical components. Like a linguistic grammar, this grammar allows for the recombination and reuse of its component “words.” One could substitute the point geometry layer `geom_point` with a density map layer `geom_density`, or the provided “classic” theme with a black-and-white theme `theme_bw` or even a custom user-created theme.<sup>5</sup>

---

<sup>5</sup>A gallery of such themes can be seen [here](#).

---

**Listing 1** An example of the grammar of graphics, in code.

---

```
library(ggplot2)
ggplot(iris, aes(x = Sepal.Length, y = Petal.Length, color = Species)) +
  geom_point() +
  geom_smooth() +
  scale_color_manual(values = c("red", "blue", "green")) +
  theme_classic()
```

---

The power of `ggplot2` lies in taking a chaotic array of plots and abstracting them into variations of a few recognizable forms. The `ggplot2` philosophy holds that every graphic can be decomposed into the grammar of graphics, no matter how many odd shapes or colors it might have. Thus, in both the “vocabulary,” or available geometries, and the “grammar,” `ggplot2` has acted as a unifying tool for statisticians and data scientists.

But as much as it is unifying and expansive, `ggplot2` is still limited. As an example, consider [force-network graphs](#), in which points are animated with positions determined by a simulation. At each step of the simulation the points’ next resulting placement is calculated as a result of the points’ previous placements, based on attraction towards linked points and repulsion from non-linked points. Points then become closer and further from each other on the graph depending on how “linked” they are, turning the graph into a sort of conceptual map. Such a graphic has been made popular in recent years by visualization tools like `d3.js`, but does not fit paradigmatically in the `ggplot2`’s grammar of graphics. The original specification of the grammar of graphics assumed static, non-animated plots; where would animated graphs like force networks fit in? What would be a force network’s scale, or the mapping between data and positional coordinates, given that nodes in a force network are updated with every “tick” of the simulation and the position of one node is determined by other nodes instead of attributes in the data?

Force networks are just one example of `ggplot2`’s general limitation. By consolidating graphics into a few recognizable forms and rules, `ggplot2` leaves out others. In some cases, like not including animated forms, this is mostly “accidental.” It doesn’t contradict the logic of `ggplot2` to extend it with the axis of time, and thus animation could very well have been left out of the `ggplot2` package simply to limit the package’s scope. Other cases, like having positions of geometries not depend on data itself but be deter-

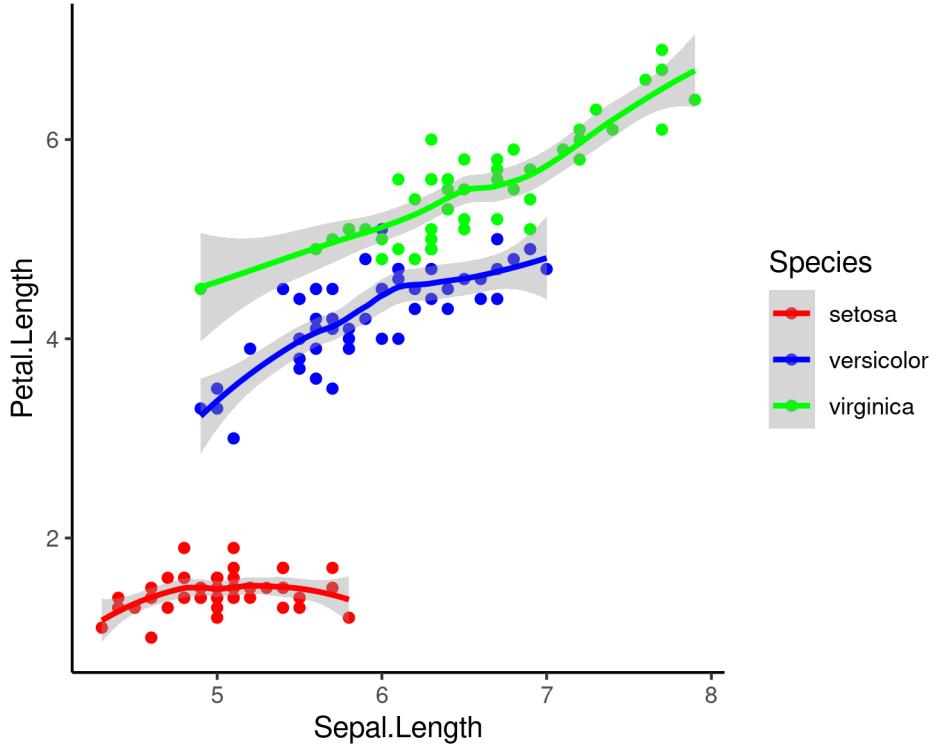


Figure 1: The rendered output of the `ggplot2` code.

mined through a random or simulated process, contradict the existing `ggplot2` assumption of having a defined scale.

The community-prescribed solution to both of these types of limitations has been to extend `ggplot2` with additional functionality and in the process rework the `ggplot2` grammar. The `ganimate` package was created to handle transitions, and is now maintained by one of the maintainers of the `ggplot2` core library.<sup>6</sup> Other libraries, like `ggigraph`, `gggraph`, `geomnet`, `ggforce`, and others have extended `ggplot2`'s functionality to visualize network graphs and other geometries that do not map perfectly to an  $x/y$  coordinate plane.

These extensions take a pragmatic position on the grammar of graphics. The grammar itself has done much work on its own to bring a coherent logic of data visualization to R, but it just as much deserves to

---

<sup>6</sup>the `ganimate` package and other animation packages in R create animations notably in a limited and performance-inefficient way compared to many tools. `ganimate` creates animations by simply appending many individual frames together, whereas web-based and OpenGL-based animation tools can use optimized rendering engines to create animations faster and with smaller file sizes. This is an issue I hope will be improved on for R in the future, but as it is not directly relevant to my project I will leave this to a footnote.

be critiqued and extended as it does to be praised. With `ggplot2` and the world of data visualization at large, relying on a few monumental perspectives have meant neglecting others.

## Du Bois

It is with this mindset that I turn to W.E.B. Du Bois, whose graphics I believe contributes many different perspectives on expressive display of information. The majority of Du Bois' work would come during the aforementioned “modern dark ages” of data visualization, during which some scholars consider few innovations to have been made. But as I will argue in this section, Du Bois's perspectives and practice of data visualization greatly preceded his time, and offers much to the data visualization world even today.

Du Bois himself was born in Massachusetts in 1868, attending an integrated public school as a child before going to Fisk University in Tennessee.<sup>7</sup> He attended Harvard University for a second degree beginning in 1888, and enrolled in graduate study at Harvard for sociology. After receiving his degree from Harvard, he began a highly prolific career with various positions at Wilberforce University, the University of Pennsylvania, Atlanta University, the Tuskegee Institute, the NAACP, and others. He died in 1963 in Ghana, while working on an encyclopedia of Africa and the African diaspora, in exile from the U.S. for his Communist sympathies.

He is known today for his massive array of contributions to the fields of African American studies and sociology, ranging from the first sociological study of a Black community in *The Philadelphia Negro*, to histories like *Africa: Its Place in Modern History*, to more theoretical texts like *The Souls of Black Folk*, and finally to creative and personal pieces like *Dusk of Dawn* and *The Quest of the Silver Fleece*. He pushed on the boundaries of all of these fields, providing new theories and questioning existing ones. This, combined with the breadth of his work in many different fields, lead some to call Du Bois “The Grandfather of Black Studies.”

Perhaps most substantially for my project, Du Bois also had unique views on statistics and quantitative information that are reflected in his work. For instance, *The Philadelphia Negro* was one of the first studies to incorporate statistics into a sociological study, now a standard practice. For data visualization,

---

<sup>7</sup>Many scholars attribute Du Bois' central interest in racism in the South to have begun during this time.

these views can be seen most clearly in Du Bois' work in the 1900 Paris Exposition, a world's fair commemorating the technological and social achievements of the 19th century. Du Bois collaborated with Daniel Murray, the Assistant Librarian of Congress at the time, and a lawyer named Thomas Calloway to bring some five hundred photographs and statistical charts in a special room titled "The Exhibit of American Negroes." Strangely alongside eugenicist "human zoos" at the Exposition that argued for the superiority of the white race, Du Bois attempted to show to some 50 million visitors the life and historical trajectory of Black people in America.<sup>8</sup>

Du Bois' graphics broke both practical and epistemological ground. Though his graphics came during the "Dark Ages of data visualization" as described by some, Du Bois still brought to the Exposition many types of shapes that had never been employed before. Like the photographs that were presented alongside them, these graphics for Du Bois represented on one hand objective events that must be seen, and an expressive and creative medium of displaying them that was just as important as any underlying factual quality. For Du Bois, it wasn't enough to convey, for example, the growth of property ownership of Black people as a simple bar chart. Du Bois saw variables like these as historical processes, where one year's data was inseparable from a previous year's, and conveyed this with the imagery of jagged spikes connecting each year's data with each other. This balance of the objective with the subjective and creative would go on to color much of Du Bois' life, leading him to comment later in works like *Dusk of Dawn* that he had embarked on a journey from "a scientist to... a master of propaganda."

There has been a substantial amount of work exploring Du Bois' work here. The most influential of these has been Whitney Battle-Baptiste and Britt Russert's *W.E.B. Du Bois' Data Portraits*, an annotated catalog of Du Bois' statistical graphics at the Exposition that has inspired a host of other works.<sup>9</sup> Anthony Starks, creator of the Go-based deck package for data visualization, began recreating Du Bois' graphics around the same time, leading many others in the data visualization community to do so as well.<sup>10</sup>.

---

<sup>8</sup>50 million visitors may sound like quite a lot, but the size and popularity of these world's fairs cannot be understated. See Rydell, *All the World's a Fair* for more information.

<sup>9</sup>Amherst, *W.E.B. Du Bois's Data Portraits*; Fusco and Olman, "Techniques of Justice"; Karduni et al., "Du Bois Wrapped Bar Chart"; Smith, "'Looking at One's Self Through the Eyes of Others"'; Haven, "W.E.B. Du Bois, Georgia, and His Data Portraits". This text, as well as the recent Artspace New Haven exhibit on Du Bois' graphics in partnership with the text's authors, served as the inspiration for this project.

<sup>10</sup>See <https://github.com/ajstarks/dubois-data-portraits> for Starks' work, Starks, "Recreating W.E.B Du Bois's Data Portraits" for his commentary, and A, "Ggplot2 Meets W. E. B. Du Bois"; Allen, "Jeremydata"; Am, Feb. 15, and 2021, "Swat

My project seeks to join this discussion with a Du Bois-inspired system of data visualization. As opposed to being a set of recreations, my project’s main contribution is a set of tools for creating graphics. These tools are reusable, in that they can be installed in the R programming language through a single line of code, and a user’s first graphics can come soon after that. They are also repurposeable and extendible, with interfaces to modify their function.

I hope to show through this project Du Bois’ choice of creative design over pure functionality, and how it might be useful for the data visualization community at large. As will be discussed further below, Du Bois often departs from the central paradigms of the grammar of graphics. He combines polar and cartesian coordinate systems together, which stretches against the notion of a single positional scale that `ggplot2` and the grammar of graphics rely on. He picks unusual axes to manipulate, like the area of a spatial feature, or new ways to manipulate existing ones, like “wrapping” a bar around itself when it reaches the end of a panel.

These often address practical problems in data visualization, like how to display items that are much larger than others. But more fundamentally, Du Bois brings a creative and open view of data visualization. Other thinkers like Edward Tufte rigidly assessed the value of a graphic by how well it approximates a single “truth,” leading to the consolidation of data visualization towards those that most clearly approximate that single truth. Systems like `ggplot2` further this, delineating clear rules and templates that lead to the canonization of certain graphics over others. Du Bois’ graphics ask us to consider an alternative, showing that there are many, many valuable ways to represent a particular set of data.

## Methods

### **ggdubois: An extension for ggplot2**

The graphics library is built using the aforementioned “grammar of graphics” implemented by the `ggplot2` package in R.

The majority of `ggdubois` contains additional geometries, or geoms, through which new shapes in Data Viz” for a few examples of explorations that followed.

the `ggplot2` framework can be created. In some cases this required only adding a preprocessing step and then modifying a combination of existing “geoms,” or geometries, in `ggplot2`. In other cases, this required reaching into the `grid` package over which `ggplot2` itself was built, and adding completely new shapes for `ggplot2`.

These geometries, or geoms, are:

1. `geom_scaledmap`, which represents a variable through the inflated or deflated area of a spatial feature on a map.
2. `geom_spike`, which represents quantities of a variable through rings that are connected by spikes.
3. `geom_spiral`, which represents quantities of a variable through a spiral shape.
4. `geom_pathspiral`, which represents quantities of a variable through a combination of a spiral and connected segments.
5. `geom_wrappedbar`, or a bar chart that wraps quantities onto the next row after they exceed a certain quantity.
6. `geom_wovenbar`, which interweaves two different bar charts on the same plane.
7. `geom_square`, which represents quantities of a variable through a filled square.

Besides new geometries, `ggdubois` contains two additional tools. The first is a theme, which can be applied to any `ggplot2` object and will replace existing the styling of that plot’s text, colors, and outline with aesthetics inspired by W. E. B. Du Bois’ work for the 1900 Paris Exposition. The second is a color palette, which is used in the theme but can also be used separately to color arbitrary graphics and plots with R.

## Data and source code

The `ggbduois` package uses and provides data from several administrative sources. These have been packaged into two datasets.

The first dataset contains time-series data at the county level for Georgia, containing decennial data from 1970 to 2000 on race, educational attainment, housing ownership, and employment. This dataset roughly mirrors much of Du Bois' own scope as seen in his graphics for the 1901 Paris Exposition. This dataset contains geographic features as well for each county or equivalents. The data here were provided by the International Public Use Microdata Services (IPUMS) project at the University of Minnesota, which assembles and transforms data to a common standard across many geographies and periods in time. The variables in this dataset originally came from the U.S. decennial Census, the TIGER/LINE program for shapefiles, and the annual American Community Survey, which are all programs administered by the United States Census Bureau.

The second dataset contains nationwide county-level data measured in 2015, and measures median household income, the unemployment rate, the child poverty rate, the population of color, the amount of particulate matter of less than  $2.5 \mu\text{m}$  in the atmosphere, the “rent burden” or the average proportion of income spent on rent, the high school graduation rate, and the Gini index for income inequality. This dataset aims to extend Du Bois’ commentary on socioeconomic inequality to the national level, with data Du Bois did not have access to in both subject and scope. More detailed descriptions of these two datasets can be found in Appendix 1. For the rest of this paper, the first dataset will be referred to as `georgia` and the second will be referred to as `demographics`.

These data as well as all source code for the `ggdubois` package can be found at <https://github.com/18kimn/ggdubois>, and the source code for this paper and the associated analyses can be found at <https://github.com/18kimn/revisiting-dubois>.

## **ggdubois**

The `ggdubois` package contains seven visualization extensions, a theme, and a color palette.

### **geom\_scaledmap**

The first custom geom that my package contributes is `geom_scaledmap`. Adapted from Plate 42 of *Data Portraits*, this geom receives one continuous variable and one categorical or discrete variable as an input,

as well as the coordinates of a geographic polygon that they are associated with. The geometry then plots the polygon once for each value of the discrete variable, scaling the polygon according to the value of the continuous variable.

An example of this is shown below. From the `georgia` dataset described above, year from 1970 to 2000 is used as a discrete variable and the proportion of high school graduates is used as a continuous variable.<sup>11</sup> Fulton County, Georgia is plotted four times, one for each specified year, and scaled to the relative size of proportion of high school graduates.

---

**Listing 2** Syntax of `geom_scaledmap()`

---

```
fulton_county <- filter(georgia, fips = "13121")
ggplot(fulton_county, aes(x = edu)) +
  geom_scaledmap() +
  facet_wrap(~year) +
  labs(
    title = "High school graduates in Fulton County, Georgia",
  )
```

---

<sup>11</sup>Note that, as is idiomatic in `ggplot2`, the information encoding geometric columns is automatically detected by `geom_scaledmap`.

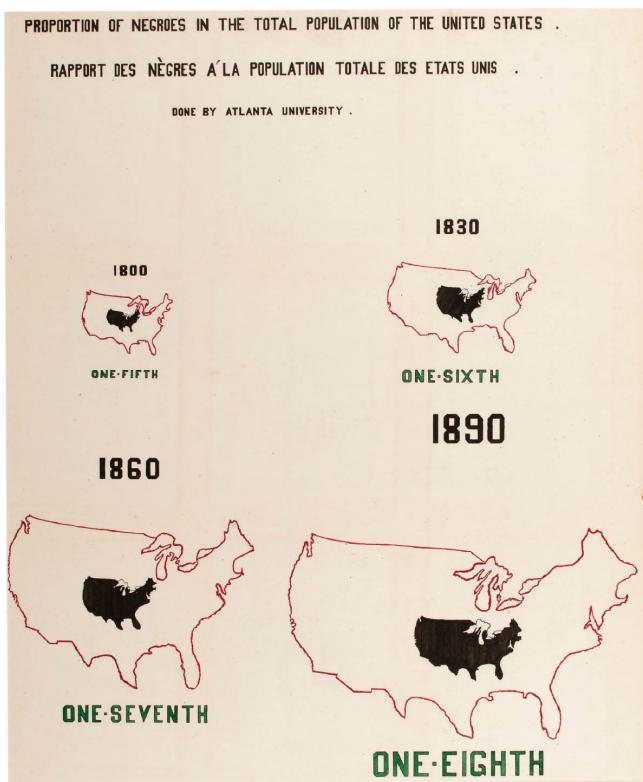
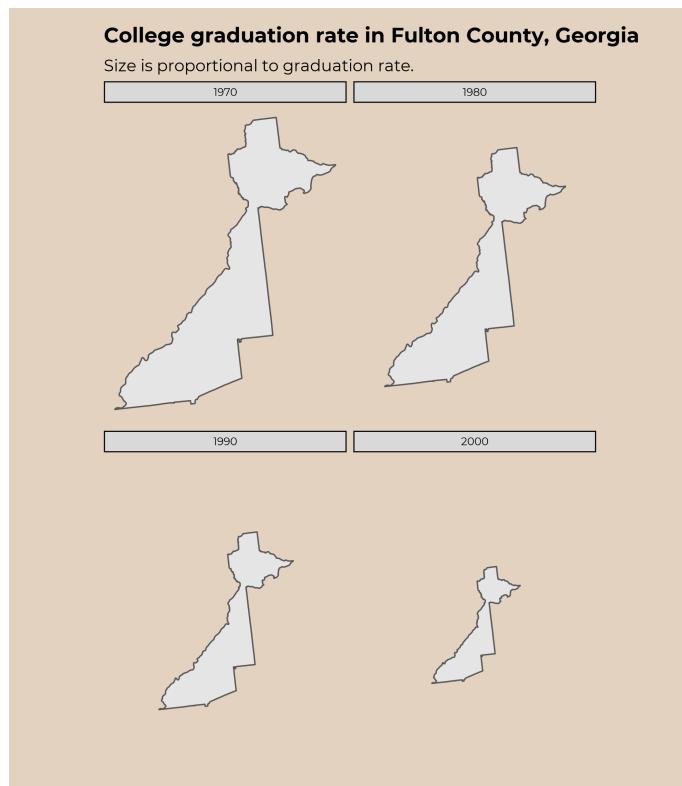


Figure 2: geom\_scaledmap

This is fairly different from a traditional map of a continuous variable, which often colors sections of a polygon instead of repeatedly drawing a polygon with modifications. Traditional maps would in this way tell viewers about the geographic distribution of a particular process, but `geom_scaledmap()` does not convey information about the geographic distribution of a process, and only simply relates a process with a geography. Viewers do not see how educational attainment in Georgia varies throughout the county, but viewers are instead presented with the association of the physical shape of Georgia with this set of measures on educational attainment.

### **geom\_spike**

The second geom that my package contributes is called `geom_spike`. This receives again one continuous variable and one categorical or discrete variable as an input. It then creates set of concurrent circles, with the distance between these concentric circles being proportional to the supplied continuous variable. “Spikes” extend from outer layers into inner ones.

---

#### **Listing 3** `geom_spike()`

---

```
ggplot(georgia, aes(x = year, y = owners)) +  
  geom_spike() +  
  coord_polar()
```

---

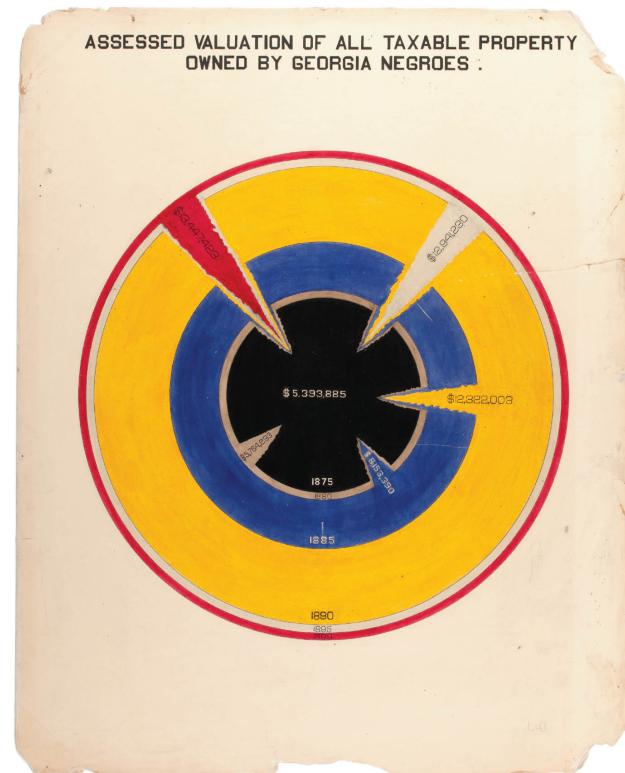
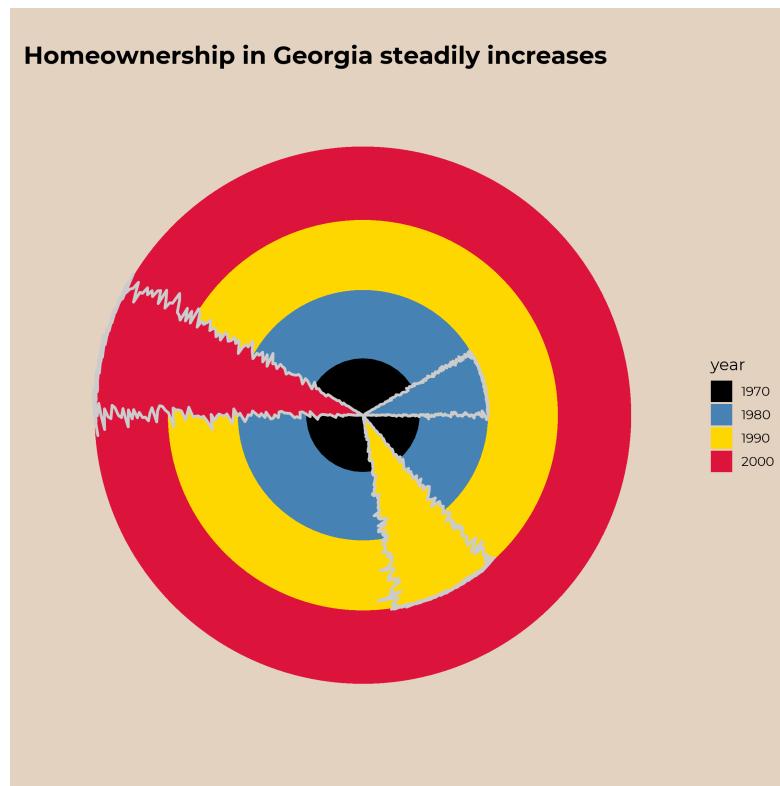


Figure 3: geom\_spike

This geom was inspired by plate 22 of *Data Portraits*, shown alongside the rendered output of `lst.`<sup>3</sup> below. Here, spikes take on a slightly more violent shape, with serrated marks stabbing into circles at the center. This graphic argues that one years' data stems from previous years' data, as opposed to being completely separated measures. In the context of my graphic, Georgia in 2000 is very much tied to Georgia in 1970.

A quirk about this geom and the `geom_spiral` function that follows is that in order to be understood within the `ggplot2` framework, the coordinate system must explicitly be specified as polar and not the default cartesian coordinate system. In the grammar of graphics, aesthetic mappings from data to coordinates and the mapping of coordinates to pixel positions are handled separately.

### **geom\_spiral**

`geom_spiral()` is inspired by Plate 25 from *Data Portraits*. This geometry is likely Du Bois' most famous, providing the title for the *Data Portraits* text and being the subject of a `#tidytuesday` tag on Twitter.<sup>12</sup>. The shape receives again one discrete variable and one continuous variable, and plots one bar for each group in a circular motion. The bars gradually curve inwards, so that a bar that wraps around the circle more than once will not intersect with itself but instead appear one level inward. The number of times that the bars will wrap can be adjusted by the user through the `nWrap` parameter supplied to `geom_spiral`.

An example of this is shown below. The `geom_spiral()` is substituted into otherwise the same syntax as the above geometries, a testament to the modular nature of `ggplot2`. The caveat to this modular nature is that to achieve the same spiral nature as Du Bois' drawings, the `coord_polar()` coordinate system must be applied. Without `coord_polar()`, a completely different (and in this case mostly nonsensical and undesired) graphic is produced.

---

<sup>12</sup>Responses and recreations to this tag can be seen at <https://github.com/rfordatascience/tidytuesday/blob/master/data/2021/2021-02-16/readme.md>

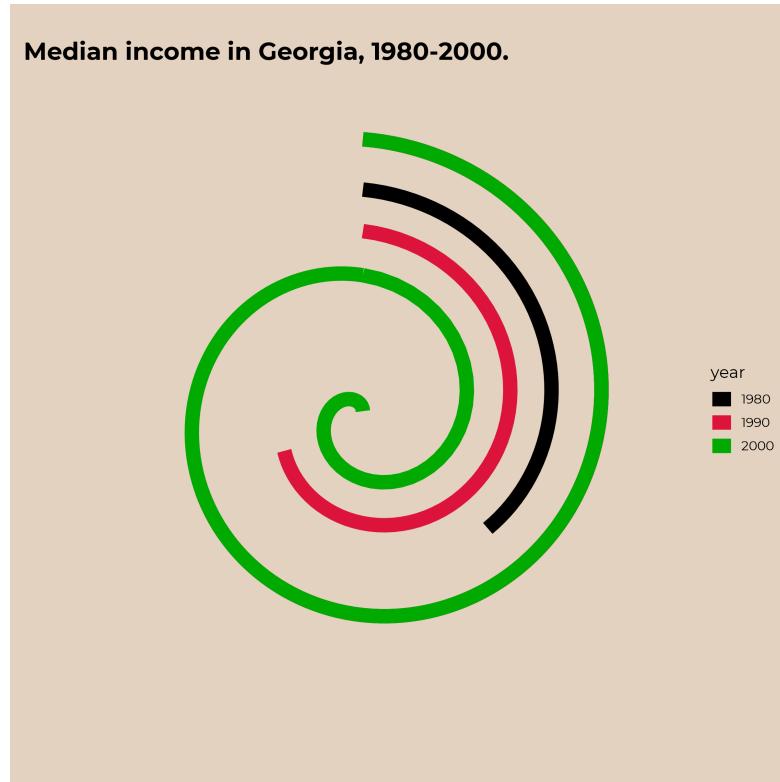
---

**Listing 4** geom\_spiral()

---

```
ggplot(georgia, aes(x = year, y = median_income)) +  
  geom_spiral() +  
  coord_polar()
```

---



The spiral shape is one of Du Bois' many responses to the issue of scale. Demographic variables like income or raw population counts usually vary wildly, often producing a graph with one or two very, very long bars and leaving other bars to be almost invisible. This would also handle space inefficiently on a panel, leaving much white space where smaller quantities are not displayed. One strategy for these cases has been the use of either focusing on only a certain range of values, or log-transforming the data before plotting. Du Bois provides an alternative strategy, making use of spirals to compactly wrap bars. With the spiral, quantities that are many times larger than others can be represented in geometrically accurate terms, meaning the area occupied by a bar is directly proportional to the quantity it represents. But unlike traditional bar plots where some bars might dwarf others, in the spiral geometry Du Bois can efficiently have longer bars take up white space that smaller bars do not use. Because of this, all quantities can usually



Figure 4: geom\_spiral

be presented alongside each other.

### **geom\_pathspiral**

The path-spiral geometry is inspired by Plate 11 from Du Bois' exhibition. Like the spiral geometry above, this shape is an answer to the problem of scale, compactly wrapping a bar so that categories associated with different values may be presented alongside each other, no matter if one category is many times larger than another. Unlike the spiral geometry shape, this shape only applies the spiral for a single category, leaving the rest of the categories to be represented in linear segments. The first linear segment is a horizontal line, and then each segment following this are connected at alternating  $+45^\circ$  and  $-45^\circ$  angles.

---

#### **Listing 5** geom\_pathspiral()

---

```
ggplot(georgia, aes(x = year, y = total_population)) +
  geom_pathspiral() +
  coord_polar()
```

---

## Population in US counties, by quartile

A work in progress



Figure 5: geom\_pathspiral

The result of this path-spiral pattern is that while the category associated with the largest value is made more compact by being wrapped around itself, the other segments take up even more room than they would on a normal coordinate plane, by stretching diagonally across the panel. Because the segments themselves do not occupy much area, Du Bois leaves much white space in between segments to add annotations like the name of a category and its associated value. My geometry does not add any annotations, but users are free to add their own through the `geom_text` and `geom_label` elements from `ggplot2`.

In this geometry Du Bois also unites polar coordinate systems and cartesian coordinate systems. The shorter line segments do not fit neatly on a polar coordinate system that the spiral geometry might imply, and the actual implementation of `geom_pathspiral` handles trigonometry under the hood so that users can avoid this issue. This geometry illustrates some creative shapes we might gain if we were to break out of the more limited coordinate system `ggplot2` allows, taking parts from different coordinate systems as it suits our purpose.

### **geom\_wrappedbar**

A simpler alternative to the question of scale is presented by Du Bois in Plates 17 and 26. While Du Bois adjusts for large quantities in some cases by wrapping bars around in a spiral fashion, as we see above, Du Bois also addresses this issue by wrapping a bar over several rows and keeping data on the cartesian plane. This has the benefit of having coordinate axes that are much more intuitively understood than a *radian/theta* combination that the `geom_spiral` geometry is plotted along. Karduni et. al have found that this strategy increases the accuracy of identifying values associated with bars, by on average 10.32 percentage points.<sup>13</sup>

The `ggdubois` implementation of this shape is in the `geom_wrappedbar` geometry. This receives again one discrete variable and one continuous variable as inputs, and creates a bar chart where all categories with values extending beyond a threshold are wrapped around the length of the panel. The default value for this threshold is one and a half times the length of the second-largest category for the width of one “cycle” of the bar, and users can supply alternative values through the optional “width” parameter.

---

<sup>13</sup>Karduni et al., “Du Bois Wrapped Bar Chart”. I believe this is also the only instance of a geometry created by Du Bois being studied in an experimental study.

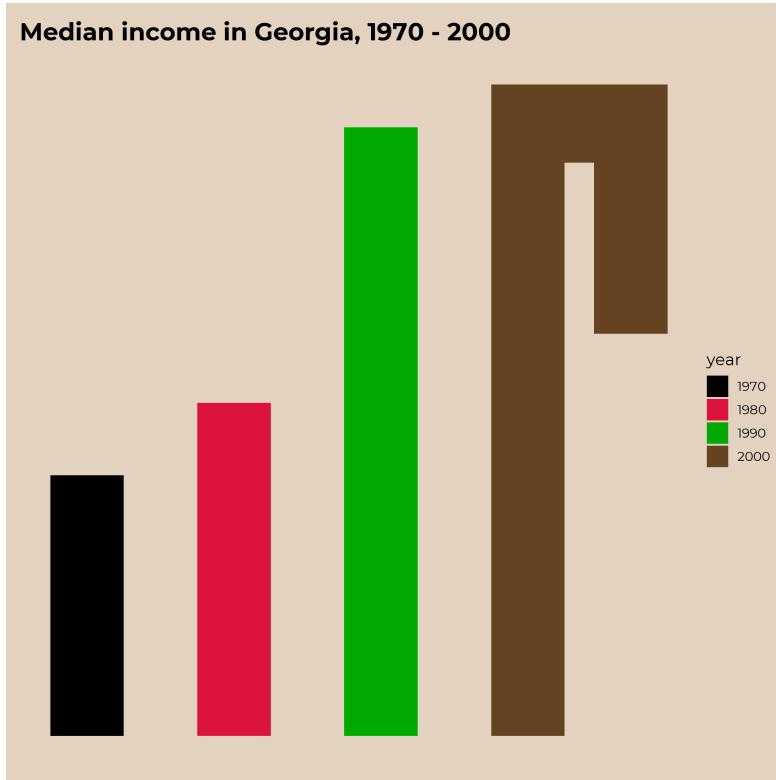
---

**Listing 6** geom\_wrappedbar()

---

```
ggplot(georgia, aes(x = year, y = median_income)) +  
  geom_wrappedbar()
```

---

**geom\_wovenbar**

geom\_wovenbar is inspired by Plate 23 of Du Bois's exhibition. It receives up to four variables: one continuous variable, up to two discrete variables with any number of levels, and one binary variable. Du Bois then interweaves two plots, one for each value of the binary variable. Each plot uses the two discrete variables as independent variables to order the x-axis, and uses the continuous variable as a response variable, drawing a bar for each possible level of the discrete variables. For a clearer explanation, see the graphic below.



Figure 6: geom\_wrappedbar

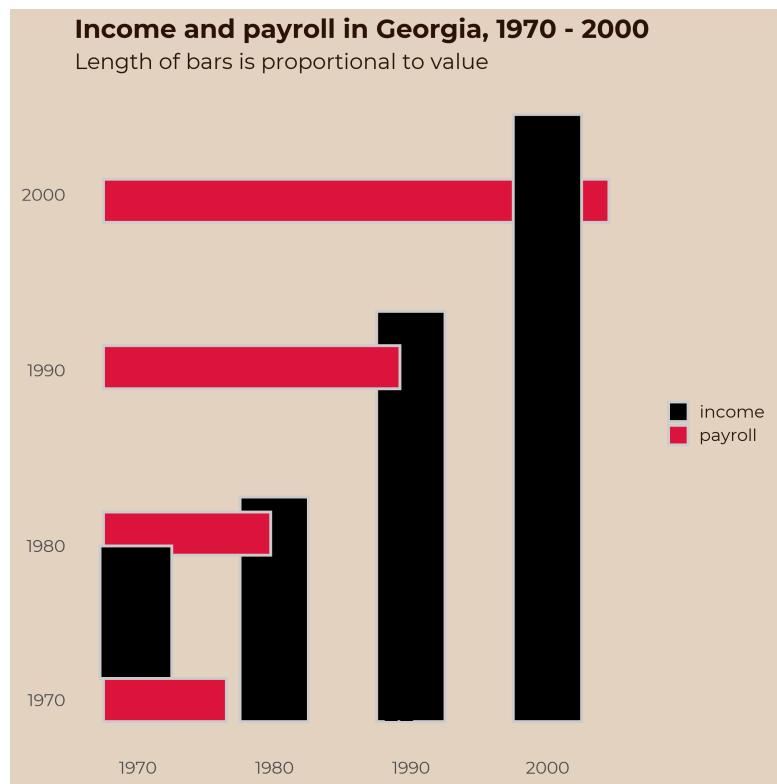
---

### Listing 7 geom\_wovenbar()

---

```
ggplot(georgia, aes(x = income, y = occupation, group = year)) +
  geom_wrappedbar()
```

---



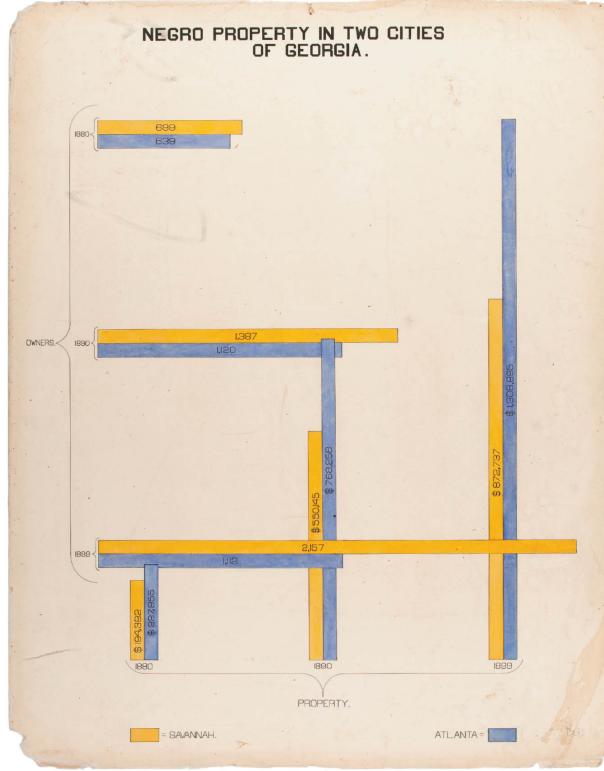


Figure 7: geom\_wovenbar

Here, the continuous variable is mapped to `y`, the binary variable is mapped to `group`, and the two discrete variables are mapped to `x` and `fill`.

Du Bois provides two insights in this shape. The first lies in incorporating four variables into a single graphic, a feat that is usually fairly difficult but is useful for seeing interactions of several variables at once. He does so by reusing directions; in my example, the x-axis is used both to represent year for the income variable and as the value of the occupation variable. The second insight Du Bois provides is the interwoven nature of these bars, that suggest an almost physical connection between the categories pictured. Like the spike plot above, Du Bois recognizes that much of the data available to him are not separate concepts but interconnected.

In return for the above two strategies, the `geom_wovenbar` graphic is possibly Du Bois' most complex shape. Especially compared to shapes natively provided by `ggplot2`, like a simple bar plot, this plot requires the reader to pause and digest each variable. This plot also does away with parts of the `ggplot2` grammar, like the association of one variable for a horizontal plane and one variable for a vertical plane

(i.e. one  $x$  variable and one  $y$  variable), here Du Bois uses both a discrete variable and continuous variable for the  $x$ -axis, and then reverses this combination for the  $y$ -axis.

### **geom\_square**

The last geometry contributed by my package is the `geom_square` function, inspired by Plate 51 of Du Bois' exhibition. It receives two continuous variables<sup>14</sup> and one categorical variable, and draws a square in which the two continuous variables order the  $x$  and  $y$  axes, and the categorical variable determines the colors and heights of different sections of the square.

An example is shown in the code snippet and graphic below. Here, the categorical variable is race, and the continuous variables are time and the percent of the total population that each race is associated with.

---

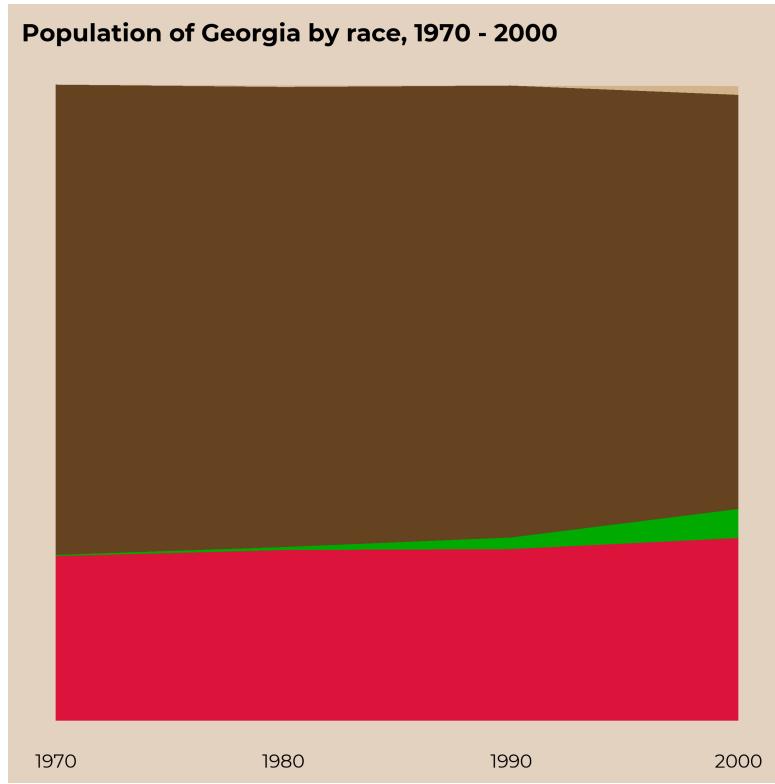
#### **Listing 8** `geom_square()`

---

```
ggplot(georgia, aes(x = year, y = percent_population, fill = race)) +  
  geom_wrappedbar()
```

---

<sup>14</sup>Both variables may also be ordered factor variables that are not truly “continuous” and only have a few values; these will be interpreted as continuous variables. Du Bois uses the ordinal variable of year in this way, whereas in other graphics year is used as a discrete variable to group separate geometries.



This graphic makes a visual argument for the tight connections between groups. For Du Bois, the population of free Black people was not separable from the population who were slaves, but tightly connected.

### **theme\_dubois**

Besides geometries, `ggdubois` also provides utilities for creating graphics in the style of W.E.B. Du Bois. The first is a theme, which can be appended to a `ggplot2` object to apply styling related to spacing between plot elements, the grid lines of the plot, the plot typography, and certain baseline colors.<sup>15</sup> Themes are overridable, meaning that users can use `theme_dubois` as a baseline theme and then tweak aesthetics as the user best sees fit.

### **dubois\_pal**

The other non-geometric utility that `ggdubois` provides is a set of color palettes: one divergent palette and one sequential palette, convenience functions for creating continuous color scales that interpolate

---

<sup>15</sup>All of the figures in this presentation have used `theme_dubois`, so I will forgo a dedicated figure.

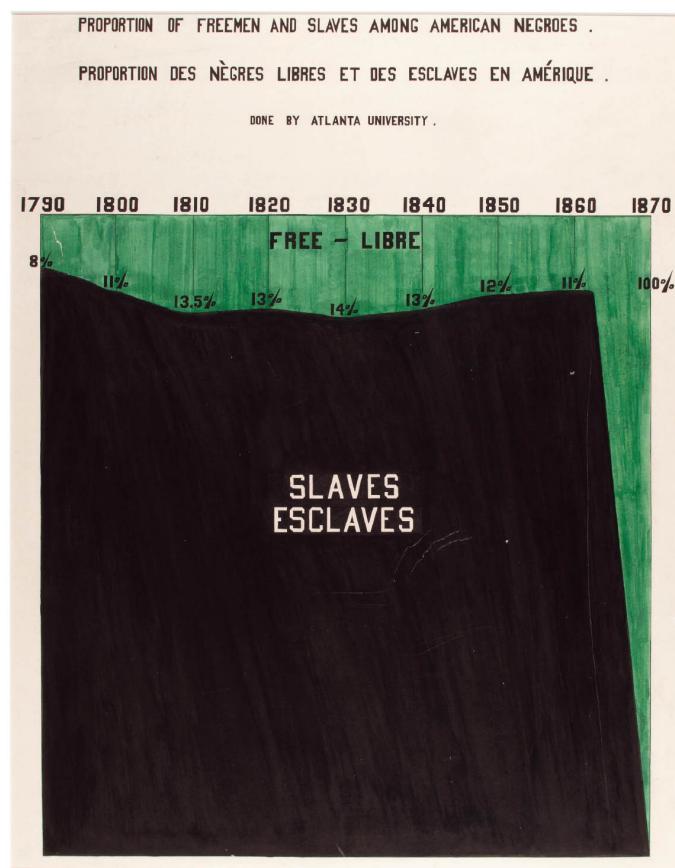


Figure 8: geom\_square

---

**Listing 9** Syntax demonstration for theme\_dubois()

---

```
ggplot(georgia, aes(x = year, y = percent_population, fill = occupation)) +  
  geom_wrappedbar() +  
  theme_dubois()
```

---

colors between the ones provided, and two functions that allow this color palette to easily be “plugged into” a `ggplot2` object.

The divergent palette is taken without modification from Anthony Starks’ Data Portraits project that reproduces many of Du Bois’ works.<sup>16</sup> Of particular note are the red, green, and black colors that are often paired together in Du Bois’ plates, which represent the colors of the Pan-African flag. Du Bois was a prominent advocate of movements like the Pan-African movement for unity against the colonial powers, related to his famous statement that the “problem of the 20th century is the problem of the color line.”

The sequential color palette is inspired by Du Bois’ colors from Plate 54, where he describes the proportion of the population that is mixed-race. For Du Bois, race is not a biological construct but an imagined one, created by political and social acts of representation. He communicates this here and tries to stretch beyond static or fixed ideas of both scales and race, using gradients at the boundaries of each polygon to symbolize how these categories are often blurry. I try to convey this perspective by providing a continuous gradient scale, as opposed to any fixed set of categories.

## Conclusion

This package still has some work required before it is ready for public use. Documentation, unit testing, additional options for each geom, and additional color palettes and styling tools are all areas for further work. Besides improving the existing components of the package, the package also has room to add more geoms by looking into Du Bois’ other work like his seminal *The Philadelphia Negro*, as well as graphics created during his time as the editor of *The Crisis* magazine.

---

<sup>16</sup>This is permissible under the Creative Commons Attribution-NonCommercial License 4.0, see <https://github.com/ajstarks/dubois-data-portraits/blob/master/LICENSE.md>.

## **``dubois\_divergent``**

A Du Bois-inspired divergent color palette.



Figure 9: The `dubois_divergent` color palette

---

**Listing 10** ggdubois color palettes.

```
# by default, returns a divergent color scale
dubois_pal(4)
# but can return a sequential color scale, e.g. with colors on a continuum
dubois_pal(10, type = "sequential")
ggplot(georgia, aes(x = median_income, y = high_school_graduates)) +
  geom_point() +
  scale_color_gradientn(colors = dubois_pal(10))
# the above is equivalent to the slightly more convenient syntax of:
ggplot(georgia, aes(x = median_income, y = high_school_graduates)) +
  geom_point() +
  scale_color_dubois()
```

---

I also hope Du Bois' perspective can be communicated in other mediums besides this package. This package is useful for creating static graphics, but just as exciting would be to see Du Bois-inspired graphics in interactive and animation-enabled mediums. The plotly library in Python and R and the Shiny framework in R are two popular examples of this, and are closely related to ggplot2 development.<sup>17</sup> The best candidate in my mind for this project is the aforementioned d3 ecosystem in JavaScript, being built to create a wildly diverse array of graphics for the web.

These additional steps can cement what this project has tried to show. As important as a grammar of graphics and a defined system of data visualization is, just as important are the creativity and tools to reach beyond this grammar. Du Bois has broken ground in the creative aspects of this endeavor; I hope through ggdubois to provide the tooling to complement this and to allow others to recognize his work.

---

<sup>17</sup>The plotly library in R can turn most graphics created with ggplot2 into an interactive graphic with a single line change, a testament to ggplot2's popularity and plotly's recognition of this. The Shiny framework for interactive data presentations is developed by RStudio, from which the ggplot2 library and many other popular systems in R also emerged

Table 1: Variables and data sources provided with the `ggdubois` package.

Variable	Source
Population counts by race and ethnicity	The IPUMS project at the University of Minnesota; the United States Census Bureau
Geographic features and shapes	IPUMS; the TIGER/LINE shapefiles program
Median household income	U.S. Census Bureau
Population of color	U.S. Census Bureau
Unemployment rate	Bureau of Labor Statistics
Child poverty rate	The County Health Rankings Project at the University of Wisconsin
Particulate matter less than 2.5 $\mu\text{m}$	The Environmental Protection Agency
Rent burden	The Eviction Lab at Princeton University

## Bibliography

- A, Matthew. "Ggplot2 Meets W. E. B. Du Bois." Stats With Matt, February 25, 2019. <https://www.statswithmatt.com/post/ggplot2-meets-w-e-b-du-bois/>.
- Allen, Jeremy. "Jeremydata: Using R's Ggplot2 to Recreate a Hand-Drawn W.E.B. Du Bois Data Visualization," February 21, 2021. <https://jeremydata.com/posts/2021-02-21-using-rs-ggplot2-to-recreate-a-hand-drawn-web-du-bois-data-visualization/>.
- Am, AUTHOR, a Luby PUBLISHED Feb. 15, and 2021. "Swat Data Viz: Du Bois Challenge." Swat Data Viz. Accessed December 17, 2021. <https://www.swarthmore.edu/NatSci/aluby1/sdv/posts/2021-02-26-dubois-challenge/>.
- Amherst, The W. E. B. Du Bois Center at the University of Massachusetts. *W. E. B. Du Bois's Data Portraits: Visualizing Black America*. Chronicle Books, 2018. <https://books.google.com?id=zft0DwAAQBAJ>.
- Bois, W. E. B. Du. *Black Folk Then and Now (The Oxford W.E.B. Du Bois): An Essay in the History and Sociology of the Negro Race*. Oxford University Press, 2014.
- . *The Philadelphia Negro: A Social Study: The Oxford W. E. B. Du Bois, Volume 2*. OUP USA, 2007.
- . *The Souls of Black Folk*. Oxford University Press, 2008.
- Bois, William Edward Burghardt Du. *Darkwater: Voices from Within the Veil*. Harcourt, Brace and Howe, 1920.
- . *The World and Africa: An Inquiry Into the Part Which Africa Has Played in World History and Color and De: The Oxford W. E. B. Du Bois*. OUP USA, 2007.
- Bois, William Edward Burghardt Du, and Henry Louis Gates. *Dusk of Dawn (the Oxford W. E. B. Du Bois)*. Oxford University Press, 2014.
- Fisher, Rebecka Rutledge. "Cultural Artifacts and the Narrative of History: W. E. B. Du Bois and the Exhibiting of Culture at the 1900 Paris Exposition Universelle." *MFS Modern Fiction Studies* 51, no. 4 (2005): 741–74. <https://doi.org/10.1353/mfs.2006.0009>.
- Friendly, Michael. "A Brief History of Data Visualization." In *Handbook of Data Visualization*,

- edited by Chun-houh Chen, Wolfgang Härdle, and Antony Unwin, 15–56. Springer Handbooks Comp.Statistics. Berlin, Heidelberg: Springer, 2008. [https://doi.org/10.1007/978-3-540-33037-0\\_2](https://doi.org/10.1007/978-3-540-33037-0_2).
- . “Milestones in the History of Data Visualization: A Case Study in Statistical Historiography.” In *Classification — the Ubiquitous Challenge*, edited by Claus Weihs and Wolfgang Gaul, 34–52. Studies in Classification, Data Analysis, and Knowledge Organization. Berlin, Heidelberg: Springer, 2005. [https://doi.org/10.1007/3-540-28084-7\\_4](https://doi.org/10.1007/3-540-28084-7_4).
- Fusco, Katherine, and Lynda C Olman. “Techniques of Justice: W. E. B. Du Bois’s Data Portraits and the Problem of Visualizing the Race.” *MELUS*, October 22, 2021, mlab031. <https://doi.org/10.1093/melus/mlab031>.
- Haven, Artspace New. “W.E.B. Du Bois, Georgia, and His Data Portraits.” Artspace. Accessed December 11, 2021. <https://artspacenewhaven.org/exhibitions/w-e-b-dubois-georgia-and-his-data-portraits/>.
- Hood, Carolyn M., Keith P. Gennuso, Geoffrey R. Swain, and Bridget B. Catlin. “County Health Rankings: Relationships Between Determinant Factors and Health Outcomes.” *American Journal of Preventive Medicine* 50, no. 2 (February 1, 2016): 129–35. <https://doi.org/10.1016/j.amepre.2015.08.024>.
- Jakubek, Joseph, and Spencer D. Wood. “Emancipatory Empiricism: The Rural Sociology of W.E.B. Du Bois.” *Sociology of Race and Ethnicity* 4, no. 1 (January 1, 2018): 14–34. <https://doi.org/10.1177/2332649217701750>.
- Karduni, Alireza, Ryan Wesslen, Isaac Cho, and Wenwen Dou. “Du Bois Wrapped Bar Chart: Visualizing Categorical Data with Disproportionate Values,” January 30, 2020. <http://arxiv.org/abs/2001.03271>.
- Lewis, David Levering, and Deborah Willis. *A Small Nation of People: W. E. B. Du Bois and African American Portraits of Progress*. Zondervan, 2010. <https://books.google.com?id=wZ9D90t7xcMC>.
- Rabaka, Reiland. “The Racialization of Information: W.E.B. Du Bois, Early Intersectionality, and Social Information.” In *Information and the History of Philosophy*. Routledge, 2021.

Rydell, Robert W. *All the World's a Fair: Visions of Empire at American International Expositions, 1876-1916*. University of Chicago Press, 2013. <https://books.google.com?id=Ab5uAAAAQBAJ>.

Shick, Greg. "The Charts of W.E.B. Du Bois - Part 1." Greg Shick, February 1, 2020. [https://gregshick.com/post/dubois-part\\_1/dubois1/](https://gregshick.com/post/dubois-part_1/dubois1/).

Smith, Shawn Michelle. "'Looking at One's Self Through the Eyes of Others': W.E.B. Du Bois's Photographs for the 1900 Paris Exposition." *African American Review* 34, no. 4 (2000): 581–99. <https://doi.org/10.2307/2901420>.

Starks, Anthony. "Recreating W.E.B Du Bois's Data Portraits." Nightingale, May 23, 2021. <https://medium.com/nightingale/recreating-w-e-b-du-boiss-data-portraits-87dd36096f34>.

Tukey, John Wilder. *Exploratory Data Analysis*. Addison Wesley Publishing Company, 1970. <https://books.google.com?id=hrJ5yAEACAAJ>.

Weheliye, Alexander G. "Diagrammatics as Physiognomy: W. E. B. Du Bois's Graphic Modernities." *CR: The New Centennial Review* 15, no. 2 (2015): 23–58. <https://doi.org/10.14321/crnewcentrevi.15.2.0023>.