

# HW #1

## Interpreting {ggplot2} code

Megan Hessel

### Table of contents

I. Setup . . . . .	2
II. Data wrangling . . . . .	3
i. Create df_pop . . . . .	3
ii. Create df_us . . . . .	3
iii. Create df_shape . . . . .	4
iv. Create df_day_hour . . . . .	5
III. Prepare text elements . . . . .	6
IV. Build plots . . . . .	6
i. Build plot_shape . . . . .	6
ii. Build plot_us . . . . .	7
iii. Build plot_day . . . . .	8
iv. Build quote*s . . . . .	9
v. Build plot_ufo . . . . .	11
vi. Build plot_base . . . . .	11
V. Assemble & save . . . . .	12
Answer some final reflective questions . . . . .	13

#### 💡 Some notes before you get started

- Be sure to install any packages in the Setup chunk that you don't already have.
- Leave the code chunk options, `eval: false` and `echo: true`, set as they are. The final infographic has been intentionally optimized (e.g., text size, spacing) for saving and viewing as a PNG file, not for display in the Plots pane or within a rendered Quarto document. As a result, the text in each individual ggplot may appear too large when viewed in the Plots pane, but will be correctly sized in the exported PNG. We'll talk more about the nuances of saving ggplots (and why these differences occur) in a later lab section.
- Questions that reference line numbers (e.g., Question #1) refer to the

line numbers shown in the rendered code chunks. You will need to render the document to view these line numbers.

- Some answers may become clearer once you've looked ahead at the code further down in the script. Consider revisiting questions as you go.

## I. Setup

```
1 library(colorspace)
2 library(geofacet)
3 library(ggtext)
4 library(ggtext)
5 library(grid)
6 library(magick)
7 library(patchwork)
8 library(scales)
9 library(showtext)
10 library(tidyverse)
11 library(glue)
12
13 ufo_sightings <- read_csv('https://raw.githubusercontent.com/rfordatascience/tidytuesday/main/2023/03/28/ufo_sightings.csv')
14 places <- read_csv('https://raw.githubusercontent.com/rfordatascience/tidytuesday/main/data/2023-03-28/locations.csv')
15
16 alien <- c("#101319", "#28ee85")
17 bg <- alien[1]
18 accent <- alien[2]
19
20 ufo_image <- magick::image_read(path = here::here("images", "ufo.png"))
21
22 sysfonts::font_add_google(name = "Orbitron", family = "orb")
23 sysfonts::font_add_google(name = "Barlow", family = "bar")
24
25 sysfonts::font_add(family = "fa-brands", regular = here::here("fonts", "Font Awesome 6 Brands"))
26 sysfonts::font_add(family = "fa-solid", regular = here::here("fonts", "Font Awesome 6 Free-Solid"))
27
28 showtext::showtext_auto(enable = TRUE)
```

1. What is the author defining in lines 15-17? Where else in the code do these defined variables show up? What advantage(s) is there to defining these values here, as variables, rather than defining the values directly throughout the script?

- *Alien* is defined as two colors which are broken down further: alien color 1 as *bg* and color 2 as *accent*. In later ggplot color parameters, these labels show up again. Defining them here is efficient. When graphing, it is easier to reference a pre-made variable label then re-writting the HEX code over and over again. Also, if the author wants to change the color later, they will only have to change it in one place.

2. **In your own words, explain what the function, `font_add_google()`, does. What's the difference between the two arguments, `name` and `family`?**

- `font_add_google` loads pre-made Google fonts into R. The `name` argument is the font name used in Google Fonts. The `family` argument is the name the author will reference later within `geom_text`.

## II. Data wrangling

### i. Create `df_pop`

```

1 df_pop <- places |>
2   filter(country_code == "US") |>
3   mutate(state = str_replace(string = state,
4                             pattern = "Fl",
5                             replacement = "FL")) |>
6   group_by(state) |>
7   summarise(pop = sum(population)) |>
8   ungroup()

```

3. **Describe what this data frame contains.**

- The `df_pop` dataframe contains 2 columns. The first column specifies each state's postal code (all capitals). The second column is each state's total population.

### ii. Create `df_us`

```

1 df_us <- ufo_sightings |>
2   filter(country_code == "US") |>
3   mutate(state = str_replace(string = state,
4                             pattern = "Fl",
5                             replacement = "FL")) |>
6   count(state) |>
7   left_join(df_pop, by = "state") |>

```

```

8   rename(num_obs = n) |>
9     mutate(
10    num_obs_per10k = num_obs / pop * 10000,
11    opacity_val = num_obs_per10k / max(num_obs_per10k)
12  )

```

#### 4. Describe what this data frame contains.

- The `df_us` dataframe has each state's total population and UFO sighting information such as total number of sighting, number of sightings per every 10k people, and the sightings in comparison to other states.

#### 5. What does `opacity_val` represent, and why is it calculated?

- `opacity_val` is a proportion of number of sightings (per every 10k people) for each state comparatively to the highest number of sightings (per every 10k people). This is calculated to have an easy state by state comparison. Just by `opacity_val`, it is easily observe Montana has the most number of sightings (per every 10k people), followed by VT, WA, WV, ME, and NH. This also sets up opacity levels (used for alpha values) representing number of sightings (per every 10k people) relative to the maximum for each state.

#### iii. Create `df_shape`

```

1 df_shape <-
2 ufo_sightings |>
3   filter(!shape %in% c("unknown", "other")) |>
4   count(shape) |>
5   rename(total_sightings = n) |>
6   arrange(desc(total_sightings)) |>
7   slice_head(n = 10) |>
8   mutate(
9     shape = fct_reorder(.f = shape,
10                         .x = total_sightings),
11     opacity_val = scales::rescale(x = total_sightings,
12                                   to = c(0.3, 1))
13   )

```

#### 6. Describe what this data frame contains.

- The `df_shape` dataframe contains the different alien sighting shapes. These shapes have the corresponding total number of sightings and the proportion of the shape's total number of sighting compared to the shape with the maximum observations.

**7. What does `fct_reorder` do when it is applied to the `shape` variable? What would have happened if this step was not performed?**

- The `fct_reorder` sorts the shape column by the number of total sightings. If this step was not performed, the dataframe would not be organized in descending order of total sighting.

**8. What is the purpose of rescaling `opacity_val`? And why rescale from 0.3 to 1?**

- The `rescale` function rescales the continuous variables from `total_sightings` between 0.3 and 1. The author choose 0.3 to 1 because they will use it later for alpha values.

**iv. Create `df_day_hour`**

```

1 df_day_hour <- ufo_sightings |>
2   mutate(
3     day = wday(reported_date_time),
4     hour = hour(reported_date_time),
5     wday = wday(reported_date_time, label = TRUE)
6   ) |>
7   count(day, wday, hour) |>
8   rename(total_daily_obs = n) |>
9   mutate(
10    opacity_val = total_daily_obs / max(total_daily_obs),
11    hour_lab = case_when(
12      hour == 0 ~ "12am",
13      hour <= 12 ~ paste0(hour, "am"),
14      hour == 12 ~ "12pm",
15      TRUE ~ paste0(hour - 12, "pm"))
16  )
17
18 colnames(df_day_hour)
```

**9. Describe what this data frame contains.**

- For every day of UFO sightning, there is sighting observations (and the proportional `opacity_val` information) for every hour.

10. What is the purpose of the last line inside the `case_when()` statement (`TRUE ~ paste0(hour - 12, "pm")`)?

- The `hour` column shows time ranging from 1-24. The `case_when` statement converts the 1-24 hour column to the typical 1-12 am and pm time slots by saying “if < 12 then its ‘am’ and if its > 12 its ‘pm’”. The `TRUE ~ paste0(hour - 12, "pm")` says “if it is anything else (aka anything greater than 12), take the hour column value (which should be between 13-24) minus it by 12 and add ‘pm’”.

### III. Prepare text elements

```
1 quotes <- paste0('"..."', str_to_sentence(ufo_sightings$summary[c(47816, 6795, 93833)]), '..."')
2
3 original <- glue("Original visualization by Dan Oehm:")
4 dan_github <- glue("<span style='font-family:fa-brands;'>&#xf09b;</span> doehm/tidytues")
5 new <- glue("Updated version by Sam Shanny-Csik for EDS 240:")
6 link <- glue("<span style='font-family:fa-solid;'>&#xf0c1;</span> eds-240-data-viz.github.io")
7 space <- glue("<span style='color:{bg};'>..</span>")
8 caption <- glue("{original}{space}{dan_github}
9             <br><br>
10            {new}{space}{link}")
```

11. In your own words, what is the difference between `paste0()` and `glue()`? Why did the author use `paste0` to construct `quotes` and `glue` to construct the other text elements?

- `paste0()` allows for the mechanical, repetitive structure within this code chunk as it creates 3 different strings of texts. `glue()` is useful because its easy to read and edit. Also, `glue()` can add formatting with HTML tags and embedded variables.

### IV. Build plots

#### i. Build `plot_shape`

```
1 plot_shape <- ggplot(data = df_shape) +
2     geom_col(aes(x = total_sightings, y = shape, alpha = opacity_val),
3               fill = accent) +
4     geom_text(aes(x = 200, y = shape, label = str_to_title(shape)),
5               family = "orb",
```

```

6         fontface = "bold",
7         color = bg,
8         size = 14,
9         hjust = 0,
10        nudge_y = 0.2) +
11      geom_text(aes(x = total_sightings - 200, y = shape, label = scales::comma(total_sightings)),
12                 family = "orb",
13                 fontface = "bold",
14                 color = bg,
15                 size = 10,
16                 hjust = 1,
17                 nudge_y = -0.2) +
18      scale_x_continuous(expand = c(0, NA)) +
19      labs(subtitle = "10 most commonly reported shapes") +
20      theme_void() +
21      theme(
22        plot.subtitle = element_text(family = "bar",
23                                      size = 40,
24                                      color = accent,
25                                      hjust = 0,
26                                      margin = margin(b = 10)),
27        legend.position = "none"
28      )

```

12. Explain the values provided to the x aesthetic for both text geoms (shape & total\_sightings).

- In the first `geom_text`, the x is 200 which describes the placement of the shape labels on the graph's y axis (200 units from the x-axis origin). In the second `geom_text`, the x is `total_sightings - 200` which describes the placement of the shape's total observations on the graph's y axis (200 units away from the end of the bar graph).

## ii. Build `plot_us`

**HINT:** Consider temporarily commenting out / rearranging the `geom_*`() layers to better understand how this plot is constructed

```

1 plot_us <- ggplot(df_us) +
2   geom_rect(aes(xmin = 0, xmax = 1, ymin = 0, ymax = 1, alpha = opacity_val),
3             fill = accent) +
4   geom_text(aes(x = 0.5, y = 0.7, label = state),
5             family = "orb",

```

```

6         fontface = "bold",
7         size = 9,
8         color = bg) +
9 geom_text(aes(x = 0.5, y = 0.3, label = round(num_obs_per10k, 1)),
10        family = "orb",
11        fontface = "bold",
12        size = 8,
13        color = bg) +
14 geofacet::facet_geo(~state) +
15 coord_fixed(ratio = 1) +
16 labs(subtitle = "Sightings per 10k population") +
17 theme_void() +
18 theme(
19   strip.text = element_blank(),
20   plot.subtitle = element_text(family = "bar",
21                               size = 40,
22                               color = accent,
23                               hjust = 1,
24                               margin = margin(b = 10)),
25   legend.position = "none"
26 )

```

13. Consider the order of `geom_*`() layers in the the above plot (`plot_us`). Why did the author order the layers in this way?

- The author builds the layers one by one. The `geom_rect()` sets up the opacity\_val alpha information. Then, with the `geom_text()`, the author builds labels: State's postal code and observation numbers. The `facet_geo()` and `coord_fix()` set up the orientation and organization of the fill and label values in the US shape. Lastly, the `labs()` and `theme()` are the finishing touches to make the visualization look cohesive.

### iii. Build `plot_day`

```

1 plot_day <- ggplot(data = df_day_hour) +
2   geom_tile(aes(x = hour, y = day, alpha = opacity_val),
3             fill = accent,
4             height = 0.9,
5             width = 0.9) +
6   geom_text(aes(x = hour, y = 9, label = hour_lab),
7             family = "orb",

```

```

8         color = accent,
9         size = 10) +
10    geom_text(aes(x = 0, y = day, label = str_sub(string = wday, start = 1, end = 1)),
11              family = "orb",
12              fontface = "bold",
13              color = bg,
14              size = 8) +
15    ylim(-5, 9) +
16    xlim(NA, 23.55) +
17    coord_polar() +
18    theme_void() +
19    theme(
20      plot.background = element_rect(fill = bg, color = bg),
21      legend.position = "none"
22    )
23
24 plot_day

```

14. This plot includes one-letter labels for each day of the week. How is this accomplished when week days are written using their three-letter abbreviations (e.g. Mon, Tue) in the df\_day\_hour data frame?

- In the label parameter, the author told R to just put the first letter of the day column with `str_sub(string = wday, start = 1, end = 1)`

15. What role do the `ylim()` and `xlim()` functions play in shaping a ggplot, and how do they change the visual layout of this particular plot? To better understand their effect, try rerunning the code with each of these lines commented out and observe how the plot's spacing and composition change.

- The `ylim()` hollows out the circle graph by creating a min and max to the y axis. The `xlim` dictate how many “slices” are in the circle graph by creating limits on the x axis.

#### iv. Build quote\*s

A comment from Dan Oehm’s original code: “A bit clunky but the path of least resistance.”

```

1 quote1 <- ggplot() +
2   annotate(geom = "text",
3           x = 0,
4           y = 1,

```

```

5   label = str_wrap(string = quotes[1], width = 40),
6   family = "bar",
7   fontface = "italic",
8   color = accent,
9   size = 16,
10  hjust = 0,
11  lineheight = 0.4) +
12  xlim(0, 1) +
13  ylim(0, 1) +
14  theme_void() +
15  coord_cartesian(clip = "off")

16
17 quote2 <- ggplot() +
18   annotate(geom = "text",
19             x = 0,
20             y = 1,
21             label = str_wrap(string = quotes[2], width = 25),
22             family = "bar",
23             fontface = "italic",
24             color = accent,
25             size = 16,
26             hjust = 0,
27             lineheight = 0.4) +
28   xlim(0, 1) +
29   ylim(0, 1) +
30   theme_void() +
31   coord_cartesian(clip = "off")

32
33 quote3 <- ggplot() +
34   annotate(geom = "text",
35             x = 0,
36             y = 1,
37             label = str_wrap(string = quotes[3], width = 25),
38             family = "bar",
39             fontface = "italic",
40             color = accent,
41             size = 16,
42             hjust = 0,
43             lineheight = 0.4) +
44   xlim(0, 1) +
45   ylim(0, 1) +
46   theme_void()

```

```
47   coord_cartesian(clip = "off")
48
49 quote1
50 quote2
```

16. Why do you think the author chose to convert these text elements (and also in `plot_ufo`, below!) into ggplot objects (you may consider returning to this question after you've worked your way through all of the code)?

- The author creates all the images and text elements prior in order to use `inset_element()` at the final plot. Breaking the visualization into various, smaller parts allows for easier manipulation.

#### v. Build `plot_ufo`

**Note:** Grob stands for graphical object. Each visual element rendered in a ggplot (e.g. lines, points, axes, entire panels, even images) is represented as a grob. Grobs can be manipulated individually to fully customize plots.

```
1 plot_ufo <- ggplot() +
2   annotation_custom(grid::rasterGrob(ufo_image)) +
3   theme_void() +
4   theme(
5     plot.background = element_rect(fill = bg, color = bg)
6   )
```

#### vi. Build `plot_base`

```
1 plot_base <- ggplot() +
2   labs(
3     title = "UFO Sightings",
4     subtitle = "Summary of over 88k reported sightings across the US",
5     caption = caption
6   ) +
7   theme_void() +
8   theme(
9     text = element_text(family = "orb",
10                       size = 48,
11                       lineheight = 0.3,
12                       color = accent),
```

```

13 plot.background = element_rect(fill = bg,
14                               color = bg),
15 plot.title = element_text(size = 128,
16                           face = "bold",
17                           hjust = 0.5,
18                           margin = margin(b = 10)),
19 plot.subtitle = element_text(family = "bar",
20                            hjust = 0.5,
21                            margin = margin(b = 20)),
22 plot.caption = ggtext::element_markdown(family = "bar",
23                                         face = "italic",
24                                         color = colorspace::darker(accent, 0.25),
25                                         hjust = 0.5,
26                                         margin = margin(t = 20)),
27 plot.margin = margin(b = 20, t = 50, r = 50, l = 50)
28 )
29
30 plot_base

```

17. Why does the author render `plot.caption` using `ggtext::element_markdown()`, rather than `element_text()` (like he does for rendering `plot.title` and `text`)?

- `element_markdown()` interprets the text as Markdown/HTMLM. Whereas `element_text()` is just for text formatting. Therefore, `element_markdown()` is better for more complex formatting, such as bolding, italics, and line breaks.

## V. Assemble & save

```

1 plot_final <- plot_base +
2   inset_element(plot_shape, left = 0, right = 1, top = 1, bottom = 0.66) +
3   inset_element(plot_us, left = 0.42, right = 1, top = 0.74, bottom = 0.33) +
4   inset_element(plot_day, left = 0, right = 0.66, top = 0.4, bottom = 0) +
5   inset_element(quote1, left = 0.5, right = 1, top = 0.8, bottom = 0.72) +
6   inset_element(quote2, left = 0, right = 1, top = 0.52, bottom = 0.4) +
7   inset_element(quote3, left = 0.7, right = 1, top = 0.2, bottom = 0) +
8   inset_element(plot_ufo, left = 0.25, right = 0.41, top = 0.23, bottom = 0.17) +
9   plot_annotation(
10     theme = theme(
11       plot.background = element_rect(fill = bg,
12                                     color = bg)
13     )

```

```

14 )
15
16 ggsave(plot = plot_final,
17     filename = here::here("outputs", "ufo_sightings_infographic.png"),
18     height = 16,
19     width = 10)

```

18. Explain how `plot_final` is assembled. What do you think is the most challenging aspect of arranging all components into a single plot?

- The `plot_final` is constructed by using `inset_element()` for all previously created plots and texts. The most challenging aspect was most likely tediously maneuvering all the pieces to fit perfectly with the parameters `left = #`, `right = #`, `top = #` `bottom = #`. Also, when examining each element, the outputs looked crazy because they were designed for the size and positions of the `plot_final`. When first creating the data visualization, it was probably difficult to size all the outputs correctly to look adequate in the final plot.

19. Can you think of one reason the author may have chosen to separate the construction of `plot_base` and `plot_final`?

- `plot_base` was used as a foundation for all the other elements to be stacked on top of in order to create the `plot_final`.

### **Answer some final reflective questions**

20. During week 2, we discuss [Choosing the right graphic form](#). Refer to this lecture when answering the sub-questions, below:

- a. What “perceptual tasks” (from Cleveland & McGill’s hierarchy) must the viewer perform to extract information from these visualizations?
  - The viewer must distinguish between shading and saturation for all the plots, especially in figure 2 and 3. Additionally, in the first figure, the viewer must also differentiate between portions.
- b. What task(s) do you think the author wanted to enable or message(s) he wanted to convey with these visualizations (see lecture 2.1, slide 16 for examples)? Be sure to note at least one task / message for each of the three data viz.

- The author is showing the audience the trends of UFO sightings - light shapes, in the early mornings, and in Virginia and Montana are the most common sighting traits. In the first figure, the author is comparing the reported shapes. In the second figure, the author is showing sightings' spatial patterns. In the third figure, the author is examining the temporal patterns of the UFO reports.
- c. Name at least one caveat to the “hierarchy of perceptual tasks” that the author employed to achieve a goal(s) you noted in question b?
- By placing the figures on a black background, it allows the viewers to differentiate between the hues/saturations easier. Also the chosen color, has high luminance when the hue is higher, letting the viewer see the trends better.
21. Describe two elements of this piece that you find visually-pleasing / easy to understand / intuitive. Why?
- I love the first bar graph because it is super clear. The observation numbers are on the graph itself and the proportions between the bars are well done. Additionally, I like the USA sightings per state figure. The figure design is visually appealing and geographically intuitive.
22. Describe two elements of this piece that you feel could be better presented in a different way. Why?
- (1) The quotes are confusing. I do not understand the quote themselves, why they are incomplete, and the author's intent behind them. (2) The plot titles need to be larger as they seem smaller than the quotes. The titles get lost in the infographic.
23. Describe two new things that you learned by interpreting / annotating this code. These could be packages, functions, or even code organizational approaches that you hadn't previously known about or considered.
- As time consuming as this was, I learned tons! Firstly, the data wrangling happening in a couple massive pipes was beautiful, concise, and easy to understand. When looking at the data visualization side, the organization was interesting - making all the elements first, then combining everything at the end. I learned the `glue()` function, and how to construct personalized themes with `theme_void() + theme()`.
24. How, if at all, did you use AI tools to help you interpret this code? Describe your approach to using these tools for this assignment. In what ways was consulting the documentation more (or less) helpful than using AI?
- My strategy for this homework was (1) comment the portion of code out and visually see the difference, (2) look up the documentation in R, (3) look up the documentation on Google, and (4) ask AI. I only used AI for the rescaling `opacity_val` question. After using all the documentation, I was still confused on why the author rescaled

0.3 to 1. Documentation for `rescale()` used nine words to losely described *what* the function did. AI could help me help me understand *why* the author choose 0.3 and 1.