

Data Visualization

Fall R '23

Dominic Bordelon, Research Data Librarian, ULS

Agenda

1. The grammar of graphics
2. Plotting distribution
3. Scatter plots
4. Comparing categories
5. Heat maps of two-way tables
6. Time series
7. Small multiples
8. Labeling, theming, and annotation

About the trainer

Dominic Bordelon, Research Data Librarian

University Library System, University of Pittsburgh

dbordelon@pitt.edu

Services for the Pitt community:

- Consultations
- Training (on-request and via public workshops)
- Talks (on-request and publicly)
- Research collaboration

Support areas and interests:

- Computer programming fundamentals, esp. for data processing and analysis
- Open Science and Data Sharing
- Data stewardship/curation
- Research methods; science and technology studies

Fall R Series

#	Date	Title
1	8/29	Getting Started with Tabular Data
2	9/5	Working with Data Frames
3	9/12	Data Visualization
4	9/19	Inference and Modeling Intro
5	9/26	Machine Learning Intro

R and RStudio Drop-In Hour

Bring your R questions!
If we don't know the answer,
we'll help you look for it.

Students, post-docs, faculty, and staff are all welcome.

For R users in the Pitt community
Tuesdays, 4:30–5:30 pm
Fall semester 2023

Hillman Library, Rm. 255

<https://bit.ly/pitt-r-office-23>

Service provided by:
University Library System,
Digital Scholarship Services



Scan QR code
for more info



 University of
Pittsburgh | Library System

DSS
DIGITAL SCHOLARSHIP SERVICES
@PITT

The ggplot2 package

“ggplot2 is a system for declaratively creating graphics, based on [The Grammar of Graphics](#). You provide the data, tell ggplot2 how to map variables to aesthetics, what graphical primitives to use, and it takes care of the details.”



...and using penguins examples

```
1 install.packages("palmerpenguins")
```

```
1 library(palmerpenguins)
```

```
2
```

```
3 # load palmerpenguins' data into your environment:
```

```
4 data(penguins)
```

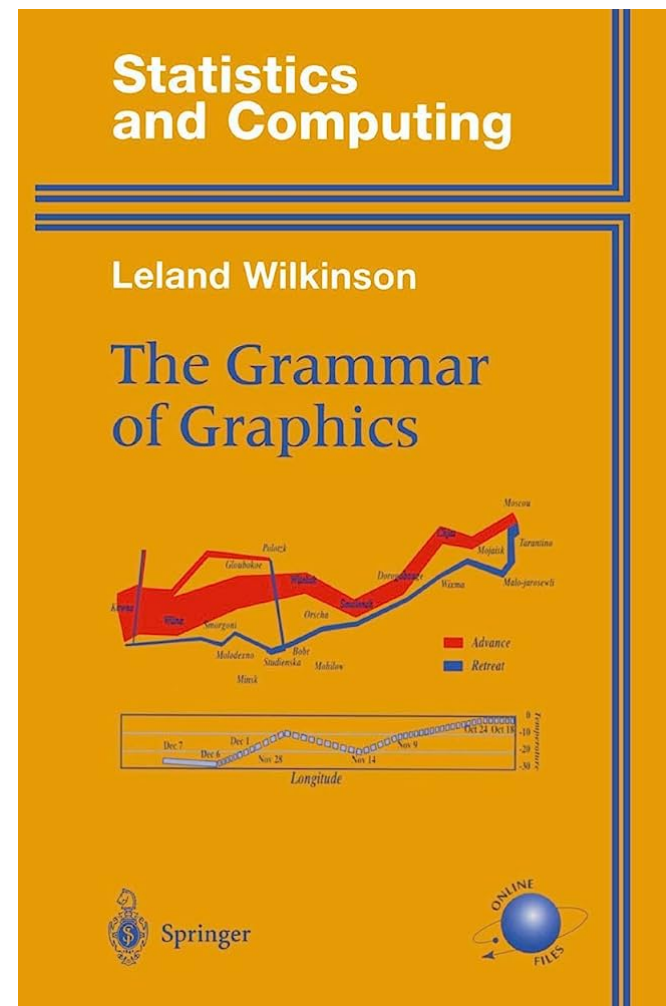
```
5 names(penguins)
```

```
[1] "species"          "island"           "bill_length_mm"  
[4] "bill_depth_mm"    "flipper_length_mm" "body_mass_g"  
[7] "sex"              "year"
```

palmerpenguins is one of many examples of an R package which functions as a downloadable data set.

The grammar of graphics

- A plot is constructed in layers:
 - data
 - aesthetics (axes, encodings)
 - scale (axis labels, color coding)
 - geometric objects (bar, scatter, heatmap tiles, etc.)
 - facets
 - statistical summaries (e.g., highlighted mean; smoother)
 - annotations
 - coordinate system (Cartesian, polar, or map projection)
 - theme



Wilkinson 1999

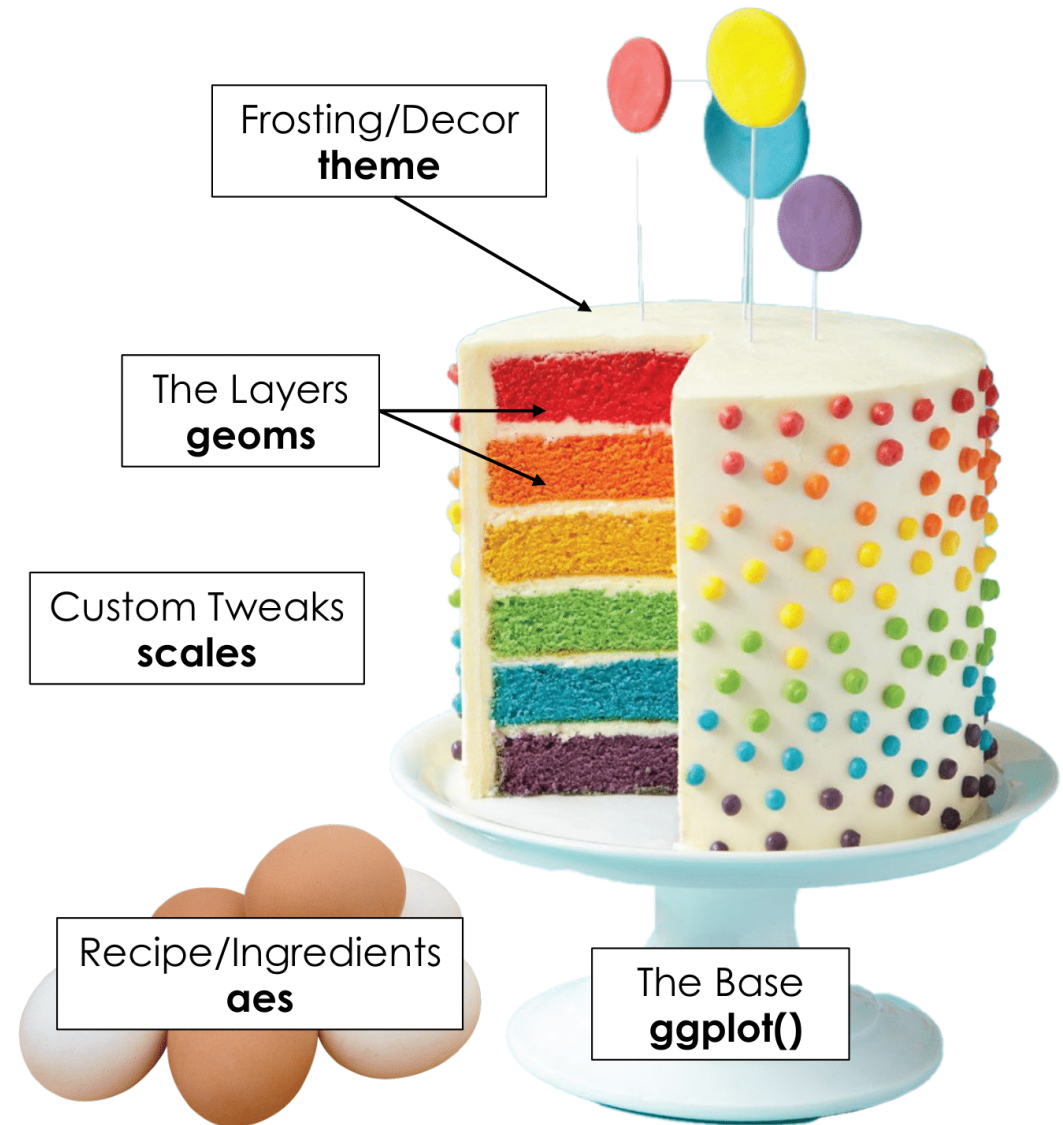
ggplot is a little bit like cake...

We *always* start by setting up the foundation with **ggplot()**

We specify our ingredients (data variables) with an **aes mapping**

We can create *layers* to our plot with **geoms**

We can style our cake ggplot with **themes**. We have out-of-the-box options, or we can go totally custom!



Source: Tanya Shapiro, https://twitter.com/tanya_shapiro/status/1576935152575340544

Aesthetic mapping

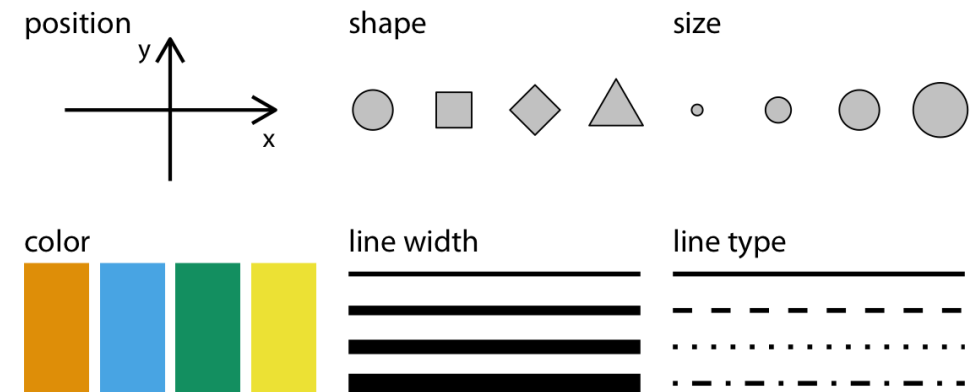
An aesthetic mapping, created with `aes()`, associates an aesthetic property of the plot with a variable in your data.

`penguins` example: I want a scatter plot of flipper length on x and bill length on y. The aesthetic mapping would be:

```
aes(x = flipper_length_mm, y =  
bill_length_mm)
```

Common aesthetics:

- `x, y`
- `color, fill`
- `linetype`
- `size`
- `shape` (of points)



Commonly used aesthetics (source: Wilke 2023)

Plotting distribution

Histograms and density plots

`geom_histogram(mapping, binwidth)` and `geom_density(mapping)`

A histogram “bins” a variable’s values and charts how many values are in each bin, giving a sense of central tendency and spread. A density plot works similarly, but it produces a smooth curve while sacrificing countable units.

```
1 # distribution of body mass:
2 penguins %>%
3   ggplot() +
4     geom_histogram(mapping = aes(x = body_mass_g))
5
6 # remove NA values, and specify bin width:
7 penguins %>%
8   drop_na(body_mass_g) %>%
9   ggplot() +
10    geom_histogram(mapping = aes(x = body_mass_g),
11                    binwidth = 100)
12
13 # density plot:
14 penguins %>%
15   ggplot() +
16     geom_density(aes(x = body_mass_g))
17
18 # layering histogram and density plot
19 # note that 1) mapping is specified in ggplot(), and inherited by geoms; 2) histogram's y mapping needs adjust
20 penguins %>%
21   drop_na(body_mass_g) %>%
22   ggplot(aes(x = body_mass_g)) +
```

Comparing distributions: violin, box

`geom_violin()`, `geom_boxplot()`. Mapping a categorical variable onto an axis will make a violin or box plot for each level of the variable.

A violin plot mirrors a density plot across the axis. A box plot marks interquartile range (white box), median (line inside the box), tails, and outliers (individual points). (where outlier has a distance $(\gt 1.5 \times \mathrm{IQR})$ from the median)

```
1 # distribution of body mass across species:
2 penguins %>%
3   ggplot() +
4     geom_violin(aes(x = body_mass_g,
5                     y = species))
6
7 penguins %>%
8   ggplot() +
9     geom_boxplot(aes(x = body_mass_g,
10                     y = species))
```

Scatter plots

geom_point() plots points

A scatter plot compares continuous variables x and y on a Cartesian plane.

`geom_point()` can take several aesthetics (`size`, `shape`, `color`), but be careful not to overload your plot with information.

```
1 penguins %>%
2   ggplot() +
3   geom_point(aes(x = flipper_length_mm,
4                 y = bill_length_mm))
5
6 # encode species as color:
7 penguins %>%
8   ggplot() +
9   geom_point(aes(x = flipper_length_mm,
10                 y = bill_length_mm,
11                 color = species))
```



geom_smooth() adds a smoother to a scatter plot

A natural next step to a scatter plot is to fit a regression model. `geom_smooth()` does this job; specify the modeling function you want with the `method` argument. The default methods are `loess`, local polynomial regression fitting (when $(n < 1000)$), or `gam`, generalized additive model with restricted maximum likelihood (when $(n \geq 1000)$). See `?stats::loess` and `?mgcv::gam` for more info on these modeling functions.

(If you have fit your own model outside of ggplot, you should instead use `predict()` to generate points to plot with `geom_line()`.)

```
1 # using default loess:
2 penguins %>%
3   ggplot(aes(x = flipper_length_mm,
4             y = bill_length_mm)) +
5   geom_point() +
6   geom_smooth()
7
8 # using a linear model, lm():
9 penguins %>%
10  ggplot(aes(x = flipper_length_mm,
11            y = bill_length_mm)) +
12  geom_point() +
13  geom_smooth(method="lm")
14
```


Comparing categories

geom_bar() makes bar charts of *counts*

```
1 penguins %>%
2   ggplot(aes(x = species)) +
3   geom_bar()
4
5 # note 2 variables, y axis, and dodge position
6 # we also reorder the factors for aesthetic reasons
7 # ...and note unexpected colors!
8 penguins %>%
9   mutate(species = fct_rev(fct_infreq(species)),
10          sex = fct_infreq(sex)) %>%
11   ggplot(aes(y = species, fill = sex)) +
12   geom_bar(position = position_dodge2(reverse=TRUE, preserve="single"))
```



Bar charts of calculated statistics

Perform summary calculations using `dplyr`, then use `geom_bar()` with `stat = "identity"` to plot the numbers as-is (default is `"count"`). ⚠ When plotting a summary statistic, one should also include error! `geom_errorbar()`, `geom_errorbarh()`, `geom_linerange()`

Example: mean body mass for each species and gender, with standard error

```
1 penguins %>%
2   drop_na(body_mass_g) %>%
3   group_by(species) %>%
4   summarize(mean_mass_g = mean(body_mass_g),
5             standard_error = sd(body_mass_g)/sqrt(n()))
6   ) %>%
7   ggplot(aes(x = species, y = mean_mass_g)) +
8   geom_bar(stat = "identity") +
9   geom_errorbar(aes(ymax = mean_mass_g + standard_error,
10                    ymin = mean_mass_g - standard_error), width = 0.5, color="blue")
11
12 # note 2 variables, y axis, and dodge position
13 # we also reorder the factors for aesthetic reasons
14 # ...and note unexpected colors!
15 penguins %>%
16   drop_na(body_mass_g) %>%
17   group_by(species, sex) %>%
18   summarize(mean_mass_g = mean(body_mass_g),
19             standard_error = sd(body_mass_g)/sqrt(n()))
20   ) %>%
21   ggplot(aes(y = species, fill = sex, x = mean_mass_g)) +
22   geom_bar(stat="identity",
```

Stacked proportional bars and pies

Sometimes we want to compare *proportions* of category membership. With `position = "fill"`, each bar has the same height, and the `fill` aesthetic works proportionately.

```
1 penguins %>%  
2   ggplot(aes(x = species, fill=sex)) +  
3   geom_bar(position = "fill") +  
4   scale_y_continuous(labels = scales::percent)
```

A pie chart in ggplot2 terms is a proportional bar cast onto a polar coordinate system.

```
1 penguins %>%  
2   mutate(species = fct_rev(fct_infreq(species))) %>%  
3   ggplot(aes(x = "", fill=species)) +  
4   geom_bar(position = "fill",  
5             color = "white",  
6             width = 1) +  
7   coord_polar("y", start=0) +  
8   theme_void()
```

Heat maps of two-way tables

Two-way tables

The *two-way table* or *contingency table* format shows the *conditional distribution* of observations among two categorical variables. In penguin terms, how many penguins of each species were observed on each of three islands?

```
1 penguins %>%  
2   janitor::tabyl(species, island)
```

species	Biscoe	Dream	Torgersen
Adelie	44	56	52
Chinstrap	0	68	0
Gentoo	124	0	0

Heat maps of two-way tables

A graphical version of the previous table is given by `geom_tile()`. However, you will need to group and summarize the data first.

```
1 penguins %>%  
2   group_by(species, island) %>%  
3   summarize(n = n()) %>%  
4   ggplot() +  
5   geom_tile(aes(y=species, x=island, fill=n))
```



Time series

geom_line() connects point observations

Time series data are the most common application for line graphs. Since line graphs expect a single y for each x , I recommend using dplyr to generate the table that you need. The examples use `n()` as the summarizing function, but one could also use `sum()` of a variable (for example).

```
1 # observations over time
2 penguins %>%
3   group_by(year) %>%
4   summarize(n = n()) %>%
5   ggplot() +
6   geom_line(aes(x=year, y=n))
7
8 # a more interesting example, using tropical storm/hurricane data
9 ?storms
10 storms %>%
11   group_by(year) %>%
12   summarize(n = n()) %>%
13   ggplot() +
14   geom_line(aes(x=year, y=n)) +
15   labs(y = "Storm observations")
```

Categories over time

Consider `geom_line()` for comparison, or `geom_area()` for showing cumulative distribution.

```
1 # observations over time
2 penguins %>%
3   group_by(year, species) %>%
4   summarize(n = n()) %>%
5   ggplot() +
6   geom_line(aes(x=year, y=n, linetype=species))
7
8 # storms by hurricane category:
9 # note: ~5-7 is the greatest number of categories you can color code before readers start having difficulty
10 storms %>%
11   mutate(category = as_factor(category)) %>%
12   drop_na(category) %>%
13   group_by(year, category) %>%
14   summarize(n = n()) %>%
15   ggplot() +
16   geom_line(aes(x=year, y=n, color=category)) +
17   labs(y = "Storm observations")
18
19 storms %>%
20   mutate(category = as_factor(category)) %>%
21   drop_na(category) %>%
22   group_by(year, category) %>%
23   summarize(n = n()) %>%
24   ggplot() +
25   geom_area(aes(x=year, y=n, fill=category)) +
26   labs(y = "Storm observations")
```

Small multiples

Small multiples

Small multiples AKA faceting replicate a plot in columns and/or rows, with one plot for each level of some categorical variable.

Want to combine multiple plots into one figure? Check out the [patchwork](#) package.

```
1 # one line graph for each storm category
2 storms %>%
3   mutate(category = as_factor(category)) %>%
4   drop_na(category) %>%
5   group_by(year, category) %>%
6   summarize(n = n()) %>%
7   ggplot() +
8   geom_line(aes(x=year, y=n, color=category)) +
9   facet_wrap(vars(category), ncol=3, nrow=2) +
10  labs(y = "Storm observations")
11
12 # a scatter for each island, with groups colored
13 penguins %>%
14   ggplot(aes(x=flipper_length_mm,
15             y=bill_length_mm,
16             color=species)) +
17   geom_point() +
18   facet_wrap(vars(island), nrow=2, ncol=2)
```

Labeling, theming, and annotation

labs() sets labels

```
1 penguins %>%  
2   drop_na(flipper_length_mm, bill_length_mm) %>%  
3   ggplot(aes(x = flipper_length_mm,  
4             y = bill_length_mm,  
5             color = species)) +  
6   geom_point() +  
7   geom_smooth(method="lm") +  
8   labs(title = "Gentoo tend to have the longest flippers",  
9        x = "Flipper length (mm)",  
10       y = "Bill length (mm)",  
11       color = "Species",  
12       caption = "Source: Gorman KB, Williams TD, Fraser WR (2014)")
```



theme_ layers change the plot's look

Start typing `theme_` in RStudio to see available options, or install and attach the `ggthemes` package for more. There is also a `theme()` function for adjusting specific aspects of the plot.

Lastly, `annotate()` can add annotations to the plot, using the coordinate system.

```
1 penguins %>%
2   drop_na(body_mass_g) %>%
3   ggplot() +
4   geom_histogram(mapping = aes(x = body_mass_g),
5                       binwidth = 100) +
6   annotate(geom="text", x = 6200, y = 7,
7               label="Biggest\n penguin") +
8   annotate(geom="segment", x = 6200, y = 5,
9               xend = max(penguins$body_mass_g, na.rm=TRUE),
10              yend = 1.5) +
11   theme_bw()
```

We can't forget to save!

`ggsave(filename, plot)` where `plot` is a `ggplot()` object you have assigned. `.png`, `.svg`, or `.pdf` in the `filename` determines the output type.

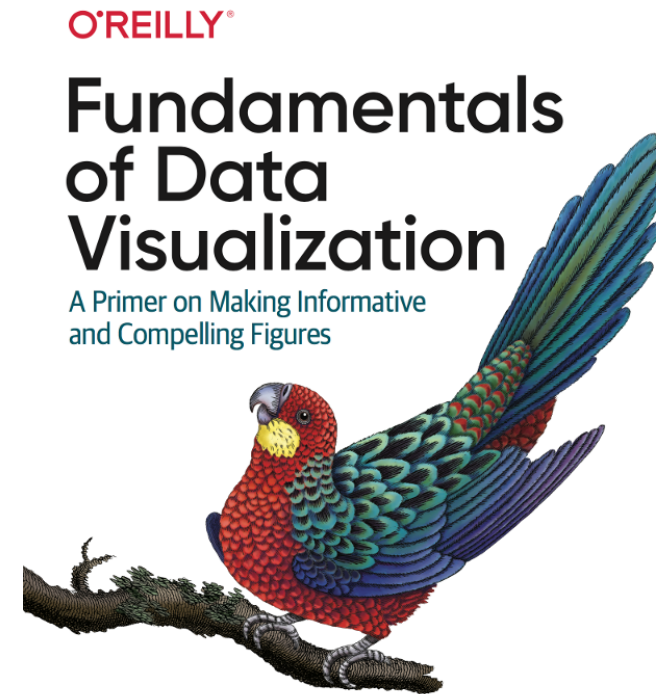
Other arguments to use: `width`, `height`, `units`

```
1 my_plot <- penguins %>%  
2   drop_na(flipper_length_mm, bill_length_mm) %>%  
3   ggplot(aes(x = flipper_length_mm,  
4             y = bill_length_mm,  
5             color = species)) +  
6   geom_point()  
7  
8 ggsave("my_plot.png", my_plot)
```



Visualize *responsibly*

- Consider the validity and reliability of your measures
- Always cite any data you visualize which you did not collect yourself
- Design for accessibility
- Study best practices, and some of the pitfalls/abuses of viz, to keep your plot honest



Claus O. Wilke

Wilke 2019, available free-to-read at
<<https://clauswilke.com/dataviz/>>

Wrap up

Session in review

Today we learned about:

- the grammar of graphics and aesthetic mapping
- lots of plot types! for a variety of applications
- some easy ways of customizing plots

Join us next week for inference and modeling!

