

Data Visualization with ggplot2

W. Evan Johnson, Ph.D.
Professor, Division of Infectious Disease
Director, Center for Data Science
Rutgers University – New Jersey Medical School

11/16/2023

Data Visualization with R

Exploratory data visualization is perhaps the greatest strength of R. One can quickly go from idea to data to plot with a unique balance of flexibility and ease.

For example, Excel may be easier than R for some plots, but it is nowhere near as flexible. D3.js may be more flexible and powerful than R, but it takes much longer to generate a plot.

Data Visualization with ggplot2

We will be creating plots using the **ggplot2**¹ package.

```
library(dplyr)  
library(ggplot2)
```

Many other approaches are available for creating plots in R. In fact, the plotting capabilities that come with a basic installation of R are already quite powerful. There are also other packages for creating graphics such as **grid** and **lattice**.

We chose to use **ggplot2** because it breaks plots into components in a way that permits beginners to create relatively complex and aesthetically pleasing plots using syntax that is intuitive and comparatively easy to remember.

¹<https://ggplot2.tidyverse.org/>

Data Visualization with ggplot2

One reason **ggplot2** is generally more intuitive for beginners is that it uses a **grammar of graphics**², the *gg* in **ggplot2**.

This is analogous to the way learning grammar can help a beginner construct hundreds of different sentences by learning just a handful of verbs, nouns and adjectives without having to memorize each specific sentence. Similarly, by learning a handful of **ggplot2** building blocks and its grammar, you will be able to create hundreds of different plots.

²<http://www.springer.com/us/book/9780387245447>

Data Visualization with **ggplot2**

Another reason **ggplot2** is easy for beginners is that its default behavior is carefully chosen to satisfy the great majority of cases and is visually pleasing. As a result, it is possible to create informative and elegant graphs with relatively simple and readable code.

One limitation is that **ggplot2** is designed to work exclusively with data tables in tidy format (where rows are observations and columns are variables).

However, a substantial percentage of datasets that beginners work with are in, or can be converted into, this format.

An advantage of this approach is that, assuming that our data is tidy, **ggplot2** simplifies plotting code and the learning of grammar for a variety of plots.

Data Visualization with ggplot2

To use **ggplot2** you will have to learn several functions and arguments. These are hard to memorize, so we highly recommend you have the ggplot2 cheat sheet handy.

You can get a copy with an internet search for “ggplot2 cheat sheet” or by clicking here:

<https://www.rstudio.com/wp-content/uploads/2015/03/ggplot2-cheatsheet.pdf>

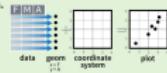
Data Visualization with ggplot2

Data Visualization with ggplot2 Cheat Sheet



Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same few components: a **data set**, a set of **geoms**—visual marks that represent data points, and a **coordinate system**.



To display data values, map variables in the data set to aesthetic properties of the geom like **size**, **color**, and **x** and **y** locations.



Build a graph with **qplot()** or **ggplot()**

qplot(x = cty, y = hwy, color = cyl, data = mpg, geom = "point")
Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

ggplot(data = mpg, aes(x = cty, y = hwy))

Begins a plot that you finish by adding layers to. No defaults, but provides more control than qplot().

data
geom of (mpg, aes(hwy, cty)) +
geom_point(aes(color = cyl)) +
geom_smooth(method = "lm") +
coord_cartesian() +
scale_color_discrete() +
theme_bw()

elements with +
add layers, default stat + layer specific aesthetic mappings
additional elements

Add a new layer to a plot with a **geom_*** or **stat_*** function. Each provides a geom, a set of aesthetic mappings, and a default stat and position adjustment.

last_plot()

Returns the last plot

ggsave("plot.png", width = 5, height = 5)

Saves last plot as 5" x 5" file named "plot.png" in working directory. Matches file type to file extension.

Geoms – Use a geom to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

One Variable

	Continuous
a	<code><- ggplot(mpg, aes(hwy))</code>
a +	<code>geom_area(stat = "bin")</code> x, y, alpha, color, fill, linetype, size b + <code>geom_area(aes(ydensity..) stat = "bin")</code> x, y, alpha, color, fill, linetype, size, weight
a +	<code>geom_density(kernel = "gaussian")</code> x, y, alpha, color, fill, linetype, size, weight b + <code>geom_density(aes(..count..))</code> x, y, alpha, color, fill
a +	<code>geom_dotplot()</code> x, y, alpha, color, fill
a +	<code>geom_freqpoly()</code> x, y, alpha, color, linetype, size b + <code>geom_freqpoly(aes(..density..))</code> x, y, alpha, color, fill, linetype, size, weight b + <code>geom_histogram(binwidth = 5)</code> x, y, alpha, color, fill, linetype, size, weight b + <code>geom_histogram(aes(..density..))</code>

	Discrete
b +	<code>geom_bar()</code> x, alpha, color, fill, linetype, size, weight
b +	<code>geom_text(label = ..label..)</code> x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

Graphical Primitives

c	<code><- ggplot(map, aes(long, lat))</code>
c +	<code>geom_polygon(aes(group = group))</code> x, y, alpha, color, fill, linetype, size

d	<code><- ggplot(economics, aes(date, unemployed))</code>
d +	<code>geom_path(linewidth = "butt", linejoin = "round", linmitre = 1)</code> x, y, alpha, color, linetype, size
d +	<code>geom_rect(pesym(ymin = unemployed - 900, ymax = unemployed + 900))</code> x, ymax, ymin, alpha, color, fill, linetype, size
d +	<code>geom_dotplot(binaxis = "y", stackdir = "center")</code> x, y, alpha, color, fill

e	<code><- ggplot(seals, aes(x = long, y = lat))</code>
e +	<code>geom_segment(aes(xend = long + delta_long, yend = lat + delta_lat))</code> x, xend, y, yend, alpha, color, linetype, size

f	<code><- geom_rect(xmin = long, ymin = lat, xmax = long + delta_long, ymax = lat + delta_lat)</code> xmax, xmin, ymin, alpha, color, fill, linetype, size
g	<code><- geom_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)</code> x, y, alpha, fill

Two Variables

Continuous X, Continuous Y

f +	<code>geom_blank()</code>
f +	<code>geom_jitter()</code> x, y, alpha, color, fill, shape, size
f +	<code>geom_point()</code> x, y, alpha, color, fill, shape, size
f +	<code>geom_quantile()</code> x, y, alpha, color, linetype, size, weight
f +	<code>geom_rug(sides = "bl")</code> alpha, color, linetype, size
f +	<code>geom_smooth(model = lm)</code> x, y, alpha, color, fill, linetype, size, weight

f +	<code>geom_text(label = ..label..)</code> x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust
-----	---

f +	<code>geom_hex()</code> x, y, alpha, colour, fill, size, weight
-----	--

Continuous Bivariate Distribution

i +	<code>geom_bin2d(binwidth = c(1, 0.5))</code> xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size, weight
i +	<code>geom_density2d()</code> x, y, alpha, colour, linetype, size
i +	<code>geom_hex()</code> x, y, alpha, colour, fill, size

Continuous Function

j +	<code>geom_area()</code> x, y, alpha, color, fill, linetype, size
j +	<code>geom_line()</code> x, y, alpha, color, linetype, size
j +	<code>geom_step(direction = "hv")</code> x, y, alpha, color, linetype, size
j +	<code>geom_crossbar(fatten = 2)</code> x, y, ymax, ymin, alpha, color, fill, linetype, size

k +	<code>geom_errorbar()</code> x, ymax, ymin, alpha, color, linetype, size, width (also <code>geom_errorbarh()</code>)
k +	<code>geom_linerange()</code> x, ymin, ymax, alpha, color, linetype, size
k +	<code>geom_pointrangle()</code> x, y, ymin, ymax, alpha, color, fill, linetype, shape, size
k +	<code>geom_rect()</code> x, y, alpha, fill, fill, fill, fill

l +	<code>geom_map(esym = TRUE, id = state, map = us.states)</code> map_id, alpha, color, fill, fill
m +	<code>geom_map(esym = TRUE, id = state, map = us.states)</code> map_id, alpha, color, fill, fill

n +	<code>geom_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)</code> x, y, alpha, fill
m +	<code>geom_tile(aes(fill = z))</code> x, y, alpha, color, fill, linetype, size

Visualizing error

df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)	<code>k < ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))</code>
--	---

o +	<code>geom_rect()</code> x, y, alpha, color, fill, fill, fill
p +	<code>geom_rect()</code> x, y, alpha, color, fill, fill, fill

q +	<code>geom_rect()</code> x, y, alpha, color, fill, fill, fill
r +	<code>geom_rect()</code> x, y, alpha, color, fill, fill, fill

s +	<code>geom_rect()</code> x, y, alpha, color, fill, fill, fill
t +	<code>geom_rect()</code> x, y, alpha, color, fill, fill, fill

u +	<code>geom_rect()</code> x, y, alpha, color, fill, fill, fill
v +	<code>geom_rect()</code> x, y, alpha, color, fill, fill, fill

w +	<code>geom_rect()</code> x, y, alpha, color, fill, fill, fill
x +	<code>geom_rect()</code> x, y, alpha, color, fill, fill, fill

y +	<code>geom_rect()</code> x, y, alpha, color, fill, fill, fill
z +	<code>geom_rect()</code> x, y, alpha, color, fill, fill, fill

Three Variables

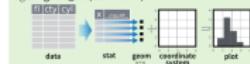
sealsSz <- with(seals, sqrt(delta_long * 2 + delta_lat * 2))	<code>m <- ggplot(seals, aes(long, lat, fill = sealsSz))</code>
m <- ggplot(seals, aes(long, lat))	<code>m + geom_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)</code> x, y, alpha, fill

m +	<code>geom_contour(aes(z = z))</code> x, y, z, alpha, colour, linetype, size, weight
m +	<code>geom_hex()</code> x, y, alpha, colour, fill, size

Data Visualization with ggplot2

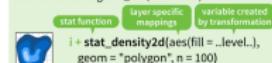
Stats - An alternative way to build a layer

Some plots visualize a transformation of the original data set. Use a stat to choose a common transformation to visualize, e.g. `a + geom_bar(stat = "bin")`



Each stat creates additional variables to map aesthetics to. These variables use a common .name.. syntax.

stat functions and geom functions both combine a stat with a geom to make a layer, i.e. `stat_bin(geom="bar")` does the same as `geom_bar(stat="bin")`.



`a + stat_bin(binwidth = 1, origin = 10)` 1D distributions

`a + stat_bindot(binwidth = 1, binsizes = "x")` x,y

`a + stat_density(binwidth = 1, kernel = "gaussian")` x,y,..count..density..scaled..

`f + stat_binned(bins = 30, drop = TRUE)` 2D distributions

`f + stat_spoke(angle = 20)` angle, radius, x, send_x, y, send_y

`f + stat_summary_hex(send_x = 30, fun = mean)` x,y,z,fill ..volume..

`f + stat_density2d(descr = 2, bins = 30, fun = mean)` x,y,fill ..volume..

`g + stat_bospital(level = 1.5)` Comparisons

`x,y | lower...middle...upper...outliers`

`g + stat_idensity(density = 1, kernel = "gaussian", scale = "area")` x,y ..density..scaled..count..n..width..width..

`f + stat_ecdf(n = 40)` Functions

`x,y | quantile..percentiles = c(0.25, 0.5, 0.75), formula = y ~ log(x), n = 40`

`f + stat_smooth(method = "auto", formula = y ~ x, se = TRUE, n = 80, fullrange = FALSE, level = 0.95)`

`x,y | ..red..x,y..ymin..ymax..`

`ggplot() + stat_function(fun = ...)` General Purpose

`fun = dnorm, n = 101, args = list sd = 0.5)`

`x | y`

`f + stat_identity()`

`geom = stat_identity(sample = 1:100), distribution = qt, distribution_params = list(q = 0.95)`

`sample | x,y | x,y | ..`

`f + stat_sum()`

`x,y | size ..size..`

`f + stat_summary(fun.data = "mean_cl_boot")`

`f + stat_unique()`

RStudio® is a trademark of RStudio, Inc. | CC-BY RStudio | info@rstudio.com | 844-448-1212 | rstudio.com

Scales

Scales control how a plot maps data values to the visual values of an aesthetic. To change the mapping, add a custom scale.

`n ~ c + geom_bar(binwidth = 1)`

`n + scale_fill_manual(values = "blue", "orange", "yellow", "green", "red")`

`scale .. to .. aesthetic .. to .. adjust .. packaged .. scale .. to .. use ..`

`scale .. to .. use .. in .. mapping ..`

`scale .. to .. use .. in .. legend ..`

`scale .. to .. use .. in .. legend ..`

range of values to include in mapping
use to use to legend
use to use in legend
use to use in legend

General Purpose scales
Use with any aesthetic:
alpha, color, fill, linetype, shape, size

`scale .. _continuous()` - map cont' values to visual values

`scale .. _discrete()` - map discrete values to visual values

`scale .. _identity()` - use data values as visual values

`scale .. _manual(values = c(...))` - map discrete values to manually chosen visual values

X and Y location scales
Use with x/y aesthetics (shown here)

`scale_x_date(labels = date_format("%m/%d/%Y"), breaks = date_breaks("2 weeks"))` - treat x values as dates. See `date_format` for label formats.

`scale_x_datetime()` - treat x values as date times. Use same arguments as `scale_x_date`.

`scale_x_log10()` - Plot x on log10 scale

`scale_x_reverse()` - Reverse direction of x axis

`scale_x_sqrt()` - Plot x on square root scale

Color and fill scales

Discrete

`n ~ b + geom_bar(binwidth = 1)`

`n + scale_fill_brewer(palette = "Blues")`

`n + scale_fill_grey(start = 0.2, end = 0.8, na.value = "#f0f0f0")`

`n + scale_fill_gradient(colors = terrain.colors(256), na.value = "#f0f0f0")`

`n + scale_fill_gradient2(low = "#ed1c24", high = "#003366", mid = "#ffffcc", space = "Lab", na.value = "#f0f0f0")`

`n + scale_fill_discrete(labels = c("A", "B", "C"))`

Continuous

`o ~ f + geom_bar(binwidth = 1)`

`o + scale_fill_brewer(palette = "Reds")`

`o + scale_fill_gradient(low = "#ffccbc", high = "#ff9800", na.value = "#f0f0f0")`

`o + scale_fill_gradient2(low = "#ff9800", high = "#ffccbc", mid = "#ffbb78", space = "Lab", na.value = "#f0f0f0")`

Shape scales

`p ~ f + geom_point(pchshape = 1)`

`p + scale_shape_manual(values = c(1:10))`

`p + scale_shape_area(max = 10, size = 10)`

Size scales

`q ~ f + geom_point(size = 10)`

`q + scale_size_area(max = 10, size = 10)`

Coordinate Systems

`r ~ c + geom_bar()`

`n ~ c + geom_bar(binwidth = 1)`

`n + scale_fill_manual(values = "blue", "orange", "yellow", "green", "red")`

`scale .. to .. aesthetic .. to .. adjust .. packaged .. scale .. to .. use ..`

`scale .. to .. use .. in .. mapping ..`

`scale .. to .. use .. in .. legend ..`

`scale .. to .. use .. in .. legend ..`

`scale .. to .. use .. in .. legend ..`

`r + coord_flip()`

`alim, ylim`

Flipped Cartesian coordinates

`r + coord_polar(theta = "x", direction = 1)`

`theta, start, direction`

Polar coordinates

`r + coord_trans(trans = "sqrt")`

`xtrans, ytrans, lims, limy`

Transformed Cartesian coordinates. Set extras and strains to the name of a window function.

`r + coord_map(projection = "ortho", orientation = c(41, -74, 0))`

`projection, orientation, alim, ylim`

Map projections from the mapproj package (mercator (default), azequalarea, lagrange, etc.)

`r + coord_map(projection = "aeqd", lon_0 = 0, lat_0 = 45)`

`projection, lon_0, lat_0, alim, ylim`

`r + coord_map(projection = "aeqd", lon_0 = 0, lat_0 = 45, r = 1000)`

`projection, lon_0, lat_0, r, alim, ylim`

`r + coord_map(projection = "aeqd", lon_0 = 0, lat_0 = 45, r = 1000, scale = 1)`

`projection, lon_0, lat_0, r, scale, alim, ylim`

`r + coord_map(projection = "aeqd", lon_0 = 0, lat_0 = 45, r = 1000, scale = 1, fill = "white")`

`projection, lon_0, lat_0, r, scale, fill, alim, ylim`

`r + coord_map(projection = "aeqd", lon_0 = 0, lat_0 = 45, r = 1000, scale = 1, fill = "white", bg = "grey")`

`projection, lon_0, lat_0, r, scale, fill, bg, alim, ylim`

`r + coord_map(projection = "aeqd", lon_0 = 0, lat_0 = 45, r = 1000, scale = 1, fill = "white", bg = "grey", fill_bg = "white")`

`projection, lon_0, lat_0, r, scale, fill, bg, fill_bg, alim, ylim`

`r + coord_map(projection = "aeqd", lon_0 = 0, lat_0 = 45, r = 1000, scale = 1, fill = "white", bg = "grey", fill_bg = "white", fill_fg = "black")`

`projection, lon_0, lat_0, r, scale, fill, bg, fill_bg, fill_fg, alim, ylim`

`r + coord_map(projection = "aeqd", lon_0 = 0, lat_0 = 45, r = 1000, scale = 1, fill = "white", bg = "grey", fill_bg = "white", fill_fg = "black", fill_fg_bg = "white")`

`projection, lon_0, lat_0, r, scale, fill, bg, fill_bg, fill_fg, fill_fg_bg, alim, ylim`

`r + coord_map(projection = "aeqd", lon_0 = 0, lat_0 = 45, r = 1000, scale = 1, fill = "white", bg = "grey", fill_bg = "white", fill_fg = "black", fill_fg_bg = "white", fill_fg_fg = "black")`

`projection, lon_0, lat_0, r, scale, fill, bg, fill_bg, fill_fg, fill_fg_bg, fill_fg_fg, alim, ylim`

`r + coord_map(projection = "aeqd", lon_0 = 0, lat_0 = 45, r = 1000, scale = 1, fill = "white", bg = "grey", fill_bg = "white", fill_fg = "black", fill_fg_bg = "white", fill_fg_fg = "black", fill_fg_fg_bg = "white")`

`projection, lon_0, lat_0, r, scale, fill, bg, fill_bg, fill_fg, fill_fg_bg, fill_fg_fg, fill_fg_fg_bg, alim, ylim`

Faceting

Facets divide a plot into subplots based on the values of one or more discrete variables.

`t ~ ggplot(mpg, aes(cty, hwy)) + geom_point()`

`t + facet_grid(~ fl)`

facet into columns based on fl

`t + facet_grid(~ year)`

facet into rows based on year

`t + facet_grid(~ fl * year)`

facet into both rows and columns

`t + facet_wrap(~ fl)`

wrap facets into a rectangular layout

Set `scales` to let axis limits vary across facets

`t + facet_grid(~ x, scales = "free")`

x and y axis limits adjust to individual facets

“free_x” - x axis limits adjust

“free_y” - y axis limits adjust

Set `labeler` to adjust facet labels

`t + facet_grid(~ fl, labeler = label_both)`

label both

`t + facet_grid(~ fl, labeler = label_bquote(alpha ^ {..x..}))`

label with alpha

`t + facet_grid(~ fl, labeler = label_parsed)`

label with parsed

Labels

`t + ggtitle("New Plot Title")`

Add a main title above the plot

`t + xlab("New X label")`

Change the label on the X axis

`t + ylab("New Y label")`

Change the label on the Y axis

`t + labs(title = "New title", x = "New x", y = "New y")`

All of the above

Legends

`t + theme(legend.position = "bottom")`

Place legend at "bottom", "top", "left", or "right"

`t + guides(color = "none")`

Set legend type for each aesthetic: colorbar, legend, or none (or legend)

`t + scale_fill_discrete(name = "Title", labels = c("A", "B", "C"))`

Set legend title and labels with a scale function.

Themes

`t + theme_bw()`

White background with gray grid lines

`t + theme_minimal()`

White background with no grid lines

`t + theme_classic()`

White background with grid lines

`t + theme_gray()`

Gray background (default theme)

`ggthemes - Package with additional ggplot2 themes`

Zooming

Without clipping (preferred)

`t + coord_cartesian(xlim = c(0, 100), ylim = c(0, 20))`

With clipping (removes unseen data points)

`t + xlim(0, 100) + ylim(0, 20)`

`t + scale_x_continuous(limits = c(0, 100)) + scale_y_continuous(limits = c(0, 100))`

Learn more at docs.ggplot2.org • ggplot2 0.9.3.1 • Updated: 3/15/2014

The Components of a Graph

We will construct the following graph for the US murders dataset:



The Components of a Graph

We can clearly see a relationship between murder totals and population size. A state falling on the dashed grey line has the same murder rate as the US average. The four geographic regions are denoted with color, which depicts how most southern states have murder rates above the average.

This data visualization shows us pretty much all the information in the data table. The code needed to make this plot is relatively simple. We will learn to create the plot part by part.

The Components of a Graph

The first step in learning ggplot2 is to be able to break a graph apart into components. The main three components to note are:

- **Data:** The US murders data table is being summarized. We refer to this as the **data** component.
- **Geometry:** The plot above is a scatterplot. This is referred to as the **geometry** component. Other possible geometries are barplot, histogram, smooth densities, qqplot, and boxplot. We will learn more about these in the Data Visualization part of the book.
- **Aesthetic mapping:** The plot uses several visual cues to represent the information provided by the dataset. The two most important cues in this plot are the point positions on the x-axis and y-axis, which represent population size and the total number of murders, respectively. Each point represents a different observation, and we *map* data about these observations to visual cues like x- and y-scale. Color is another visual cue that we map to region. We refer to this as the

The Components of a Graph

We also note that:

- The points are labeled with the state abbreviations.
- The range of the x-axis and y-axis appears to be defined by the range of the data. They are both on log-scales.
- There are labels, a title, a legend, and we use the style of The Economist magazine.

We will now construct the plot piece by piece.

The Components of a Graph

We start by loading the dataset:

```
library(dslabs)  
data(murders)
```

ggplot objects

The first step in creating a **ggplot2** graph is to define a ggplot object. We do this with the function `ggplot`, which initializes the graph. If we read the help file for this function, we see that the first argument is used to specify what data is associated with this object:

```
ggplot(data = murders)
```

We can also pipe the data in as the first argument. So this line of code is equivalent to the one above:

```
murders %>% ggplot()
```

Geometries

In ggplot2 we create graphs by adding *layers*. Layers can define geometries, compute summary statistics, define what scales to use, or even change styles. To add layers, we use the symbol `+`. In general, a line of code will look like this:

`DATA %>% ggplot() + LAYER 1 + LAYER 2 + ... +
LAYER N`

Usually, the first added layer defines the geometry. We want to make a scatterplot. What geometry do we use?

Taking a quick look at the cheat sheet, we see that the function used to create plots with this geometry is `geom_point`.



Aesthetic mappings

Aesthetic mappings describe how properties of the data connect with features of the graph, such as distance along an axis, size, or color. The aes function connects data with what we see on the graph by defining aesthetic mappings and will be one of the functions you use most often when plotting. The outcome of the aes function is often used as the argument of a geometry function. This example produces a scatterplot of total murders versus population in millions:

```
murders %>% ggplot() +  
  geom_point(aes(x = population/10^6, y = total))
```

We can drop the x = and y = if we wanted to since these are the first and second expected arguments, as seen in the help page.

Instead of defining our plot from scratch, we can also add a layer to the p object that was defined above as p <- ggplot(data = murders):

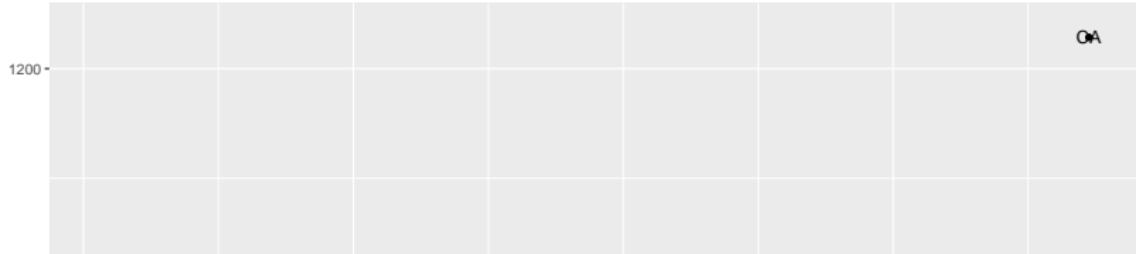
```
p + geom_point(aes(population/10^6, total))
```

Layers

A second layer in the plot we wish to make involves adding a label to each point to identify the state. The `geom_label` and `geom_text` functions permit us to add text to the plot with and without a rectangle behind the text, respectively.

Because each point (each state in this case) has a label, we need an aesthetic mapping to make the connection between points and labels. By reading the help file, we learn that we supply the mapping between point and label through the `label` argument of `aes`. So the code looks like this:

```
p + geom_point(aes(population/10^6, total)) +
  geom_text(aes(population/10^6, total, label = abb))
```



Global versus local aesthetic mappings

In the previous line of code, we define the mapping

`aes(population/10^6, total)` twice, once in each geometry.

We can avoid this by using a *global* aesthetic mapping. We can do this when we define the blank slate `ggplot` object. Remember that the function `ggplot` contains an argument that permits us to define aesthetic mappings:

```
args(ggplot)
```

```
## function (data = NULL, mapping = aes(), ..., environment =  
## NULL
```

If we define a mapping in `ggplot`, all the geometries that are added as layers will default to this mapping. We redefine `p`:

```
p <- murders %>% ggplot(aes(population/10^6, total, label =
```

and then we can simply write the following code to produce the previous plot:

Scales

First, our desired scales are in log-scale. This is not the default, so this change needs to be added through a *scales* layer. A quick look at the cheat sheet reveals the `scale_x_continuous` function lets us control the behavior of scales. We use them like this:

```
p + geom_point(size = 3) +
  geom_text(nudge_x = 0.05) +
  scale_x_continuous(trans = "log10") +
  scale_y_continuous(trans = "log10")
```

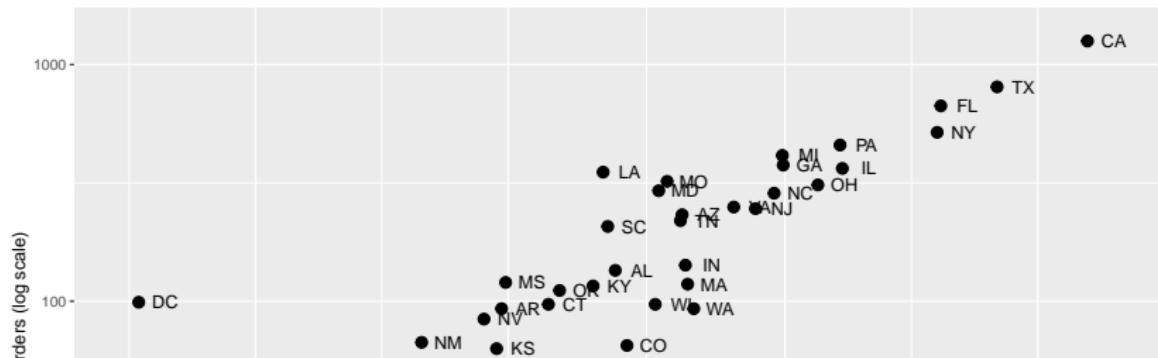


Labels and titles

Similarly, the cheat sheet quickly reveals that to change labels and add a title, we use the following functions:

```
p + geom_point(size = 3) +
  geom_text(nudge_x = 0.05) +
  scale_x_log10() +
  scale_y_log10() +
  xlab("Populations in millions (log scale)") +
  ylab("Total number of murders (log scale)") +
  ggtitle("US Gun Murders in 2010")
```

US Gun Murders in 2010



Categories as colors

We can change the color of the points using the `col` argument in the `geom_point` function. To facilitate demonstration of new features, we will redefine `p` to be everything except the points layer:

```
p <- murders %>% ggplot(aes(population/10^6, total, label  
  geom_text(nudge_x = 0.05) +  
  scale_x_log10() +  
  scale_y_log10() +  
  xlab("Populations in millions (log scale)") +  
  ylab("Total number of murders (log scale)") +  
  ggtitle("US Gun Murders in 2010")
```

and then test out what happens by adding different calls to `geom_point`. We can make all the points blue by adding the `color` argument:

```
p + geom_point(size = 3, color = "blue")
```

US Gun Murders in 2010

Annotation, shapes, and adjustments

We often want to add shapes or annotation to figures that are not derived directly from the aesthetic mapping; examples include labels, boxes, shaded areas, and lines.

Here we want to add a line that represents the average murder rate for the entire country. Once we determine the per million rate to be r , this line is defined by the formula: $y = rx$, with y and x our axes: total murders and population in millions, respectively. In the log-scale this line turns into: $\log(y) = \log(r) + \log(x)$. So in our plot it's a line with slope 1 and intercept $\log(r)$. To compute this value, we use our **dplyr** skills:

```
r <- murders %>%  
  summarize(rate = sum(total) / sum(population) * 10^6) %>%  
  pull(rate)
```

To add a line we use the `geom_abline` function. `ggplot2` uses `ab` in the name to remind us we are supplying the intercept (`a`) and slope (`b`). The default line has slope 1 and intercept 0 so we only

Add-on packages

The power of **ggplot2** is augmented further due to the availability of add-on packages. The remaining changes needed to put the finishing touches on our plot require the **ggthemes** and **ggrepel** packages.

The style of a **ggplot2** graph can be changed using the theme functions. Several themes are included as part of the **ggplot2** package. In fact, for most of the plots in this book, we use a function in the **dslabs** package that automatically sets a default theme:

```
ds_theme_set()
```

Many other themes are added by the package **ggthemes**. Among those are the `theme_economist` theme that we used. After installing the package, you can change the style by adding a layer like this:

```
library(ggthemes)
+ theme_economist()
```

Putting it all together

Now that we are done testing, we can write one piece of code that produces our desired plot from scratch.

```
library(ggthemes)
library(ggrepel)

r <- murders %>%
  summarize(rate = sum(total) / sum(population) * 10^6) %>
  pull(rate)

murders %>% ggplot(aes(population/10^6, total, label = abbrev))
  geom_abline(intercept = log10(r), lty = 2, color = "darkgreen")
  geom_point(aes(col=region), size = 3) +
  geom_text_repel() +
  scale_x_log10() +
  scale_y_log10() +
  xlab("Populations in millions (log scale)") +
  ylab("Total number of murders (log scale)") +
```

Quick plots with qplot

We have learned the powerful approach to generating visualization with ggplot. However, there are instances in which all we want is to make a quick plot of, for example, a histogram of the values in a vector, a scatterplot of the values in two vectors, or a boxplot using categorical and numeric vectors. We demonstrated how to generate these plots with `hist`, `plot`, and `boxplot`. However, if we want to keep consistent with the ggplot style, we can use the function `qplot`.

If we have values in two vectors, say:

```
data(murders)
x <- log10(murders$population)
y <- murders$total
```

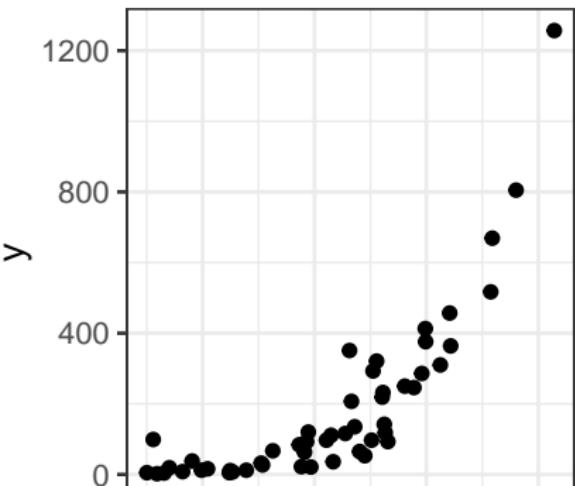
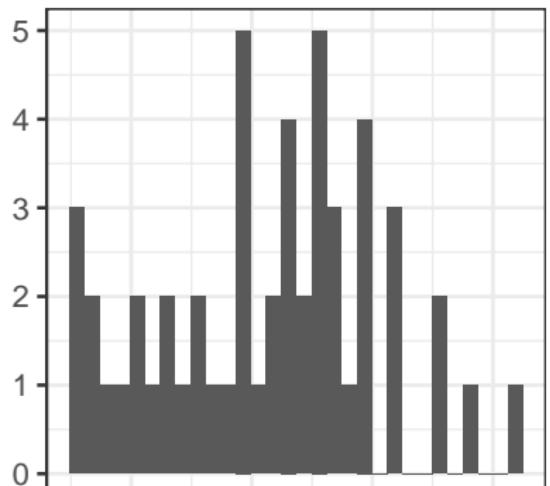
and we want to make a scatterplot with ggplot, we would have to type something like:

```
data.frame(x = x, y = y) %>%
  ggplot(aes(x, y)) +
```

Grids of plots

There are often reasons to graph plots next to each other. The **gridExtra** package permits us to do that:

```
library(gridExtra)
p1 <- qplot(x)
p2 <- qplot(x,y)
grid.arrange(p1, p2, ncol = 2)
```



Exercises

Start by loading the **dplyr** and **ggplot2** library as well as the murders and heights data.

```
library(dplyr)
library(ggplot2)
library(dslabs)
data(heights)
data(murders)
```

1. With **ggplot2** plots can be saved as objects. For example we can associate a dataset with a plot object like this

```
p <- ggplot(data = murders)
```

Because data is the first argument we don't need to spell it out

```
p <- ggplot(murders)
```

and we can also use the pipe:

```
p <- murders %>% ggplot()
```

Session Info

```
sessionInfo()
```

```
## R version 4.3.2 (2023-10-31)
## Platform: aarch64-apple-darwin20 (64-bit)
## Running under: macOS Ventura 13.5.1
##
## Matrix products: default
## BLAS: /Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/4.3-arm64/Resources/lib/libRlapack.dylib; LAPACK ver
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## time zone: America/Denver
## tzcode source: internal
##
## attached base packages:
## [1] stats      graphics   grDevices  utils       datasets   methods    base
##
## other attached packages:
## [1] gridExtra_2.3   ggrepel_0.9.4   ggthemes_4.2.4   dslabs_0.7.6
## [5] lubridate_1.9.3forcats_1.0.0  stringr_1.5.0   dplyr_1.1.3
## [9] purrr_1.0.2    readr_2.1.4    tidyr_1.3.0    tibble_3.2.1
## [13] ggplot2_3.4.4  tidyverse_2.0.0
##
## loaded via a namespace (and not attached):
## [1] gtable_0.3.4     compiler_4.3.2    Rcpp_1.0.11      tidyselect_1.2.0
## [5] scales_1.2.1     yaml_2.3.7      fastmap_1.1.1   R6_2.5.1
## [9] labeling_0.4.3   generics_0.1.3   knitr_1.45     munsell_0.5.0
## [13] pillar_1.9.0     tzdb_0.4.0      rlang_1.1.2     utf8_1.2.4
## [17] stringi_1.8.1   xfun_0.41      timechange_0.2.0 cli_3.6.1
## [21] withr_2.5.2     magrittr_2.0.3   digest_0.6.33   grid_4.3.2
## [25] rstudioapi_0.15.0hms_1.1.3    lifecycle_1.0.4 vctrs_0.6.4
## [29] evaluate_0.23   glue_1.6.2     farver_2.1.1   fansi_1.0.5
## [33] colorspace_2.1-0rmarkdown_2.25 tools_4.3.2     pkgconfig_2.0.3
```