

Bespoke Data Visualization using R and ggplot2

A CHI 2022 Course

Sandy Gould
Cardiff University
goulds@cardiff.ac.uk

New Orleans, 4th May 2022

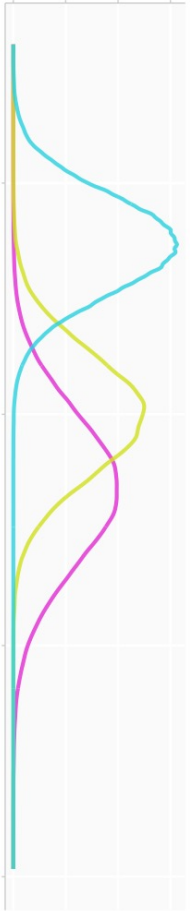
This course

Welcome!

Thank you for joining this course!

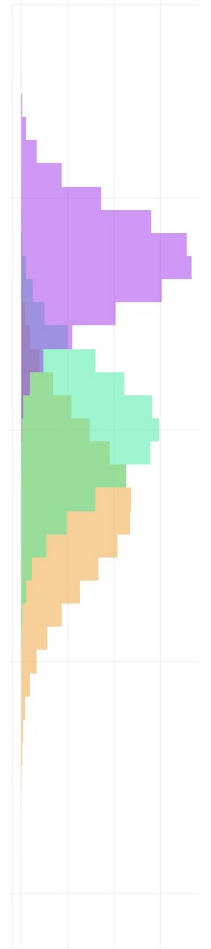
There are two parts to it:

- I want to talk to you about our approach to visualization in this course.
- Then I want us to work through a set of practical tasks. This will be the vast majority of our time.



This course

- Is an updated version of one I ran in 2019.
- It was well received in 2019.
- I've corrected some small errors in this one.
- I've made it two units; folks felt rushed last time.



Glossary

Term	Meaning
R	A statistical programming language
RStudio	An environment for writing R programs
ggplot2	A library for creating visualizations using R
Declarative visualization	A way of specifying visualizations

Materials

<http://sjjg.uk/chi22-course>

Click the link for materials.

Declarative visualizations

Specifying visualizations declaratively means that:

- Output is **predictable**
 - *e.g., if we want to redraw*
- It's easy to completely change the visualization without having to mess around with the underlying data.
- We can create more complex visualizations with a higher degree of **reliability**
- Text-based declarations will always be **readable**

The ggplot2 approach

ggplot2 takes a declarative approach to the creation of visualizations.

The emphasis is on:

- A logical structure to the creation of visualizations
- Structuring of the data that is being used
- High quality defaults

'The Grammar of Graphics'

Leland Wilkinson published 'The Grammar of Graphics', an influential text on the construction of graphics. Hadley Wickham, the developer of of ggplot2, has taken Wilkinson's philosophy into ggplot2.

Essentially, everything is split into:

- **Aesthetics** (*what is going to be shown and how it will be shown*)
- **Geometries** (*the form of plots that will represent the data*)

I find the labels for these confusing at times!

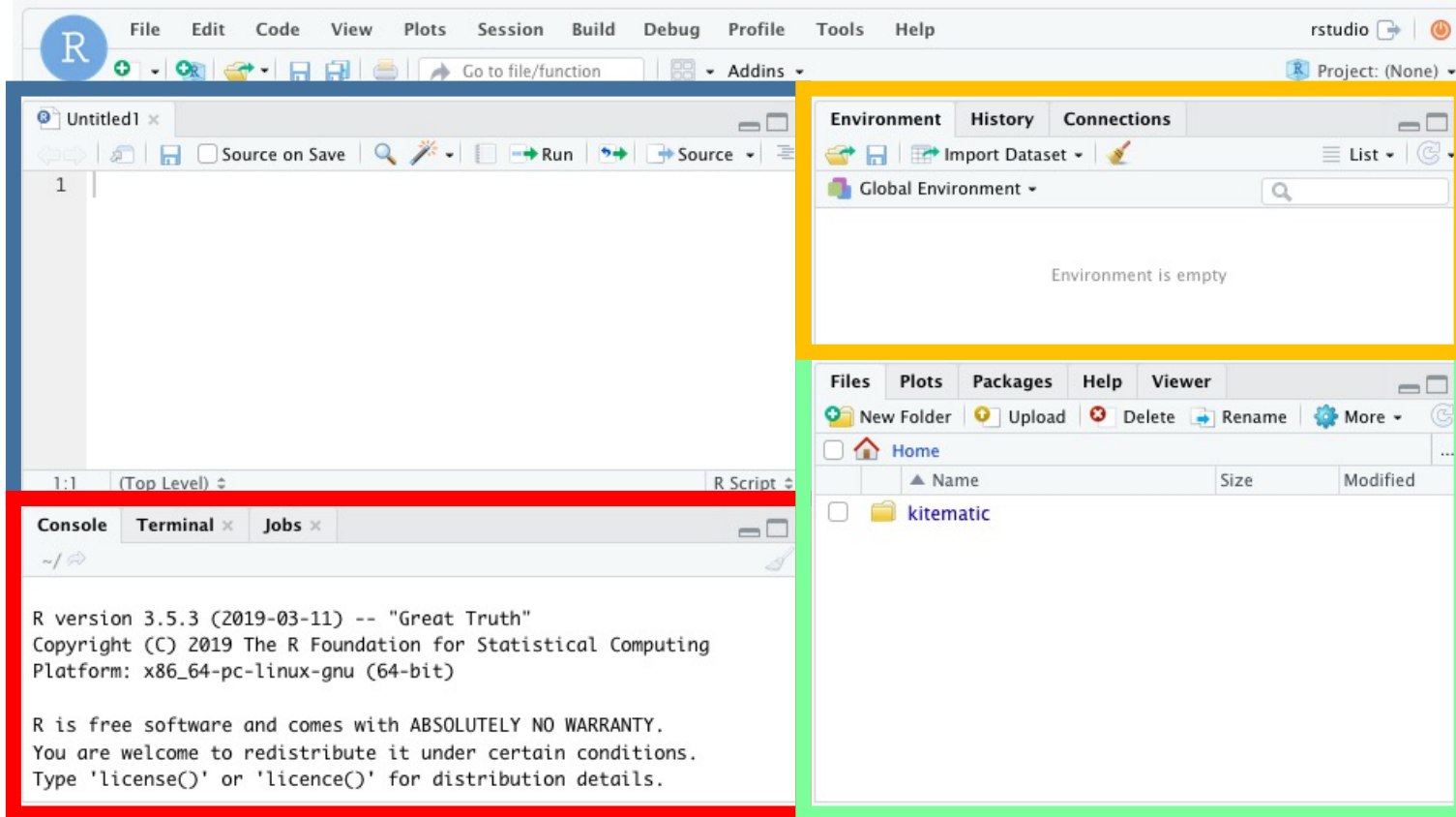
Writing code for R

- R is written using plaintext files.
- They can be written using a text editor and run from the command line.
- They can be written and run from R's default editing environment.
- They can be written and run from RStudio, a complete development environment for R.

Writing code for R

- The goal of today is not to teach you programming.
- It is to give you an understanding of how plots can be pieced together using a variety of ggplot2 commands.
- You will leave with an understanding of the logic of ggplot2 as well as plenty of template code to get you started.
- The website will stay up for you to refer to in future.

RStudio



Getting RStudio up and running

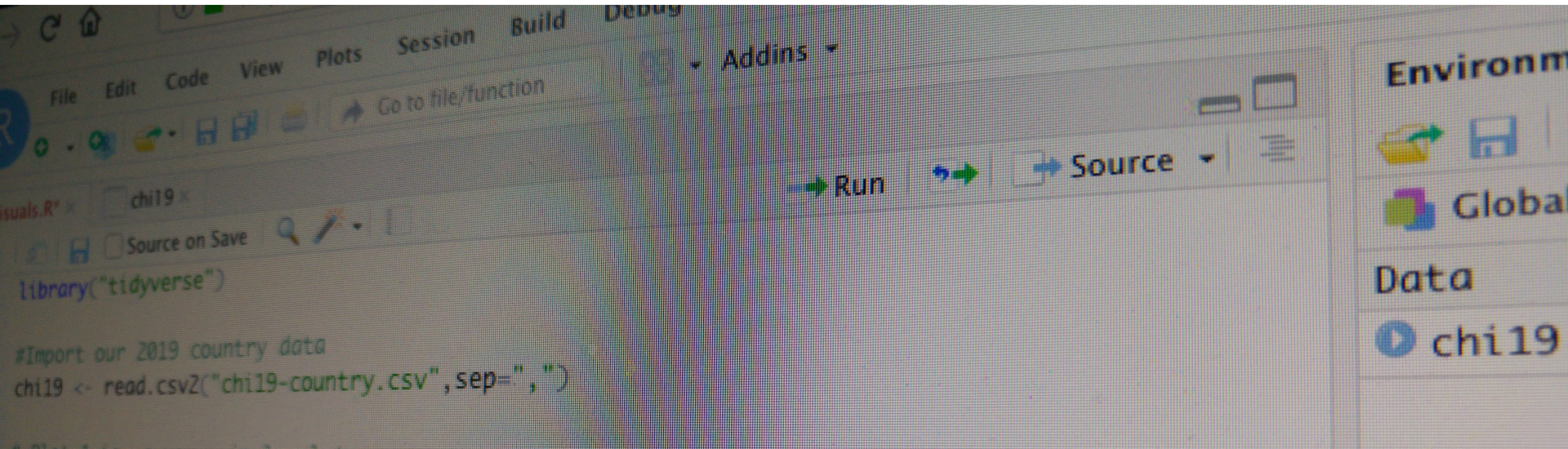
We're going to be using RStudio in the browser.

<https://rstudio.sjg.uk>

Username: friend+<your number> (e.g., friend1, friend2)

Password: ggplot2

Let's have a look at an example together



Our dataset

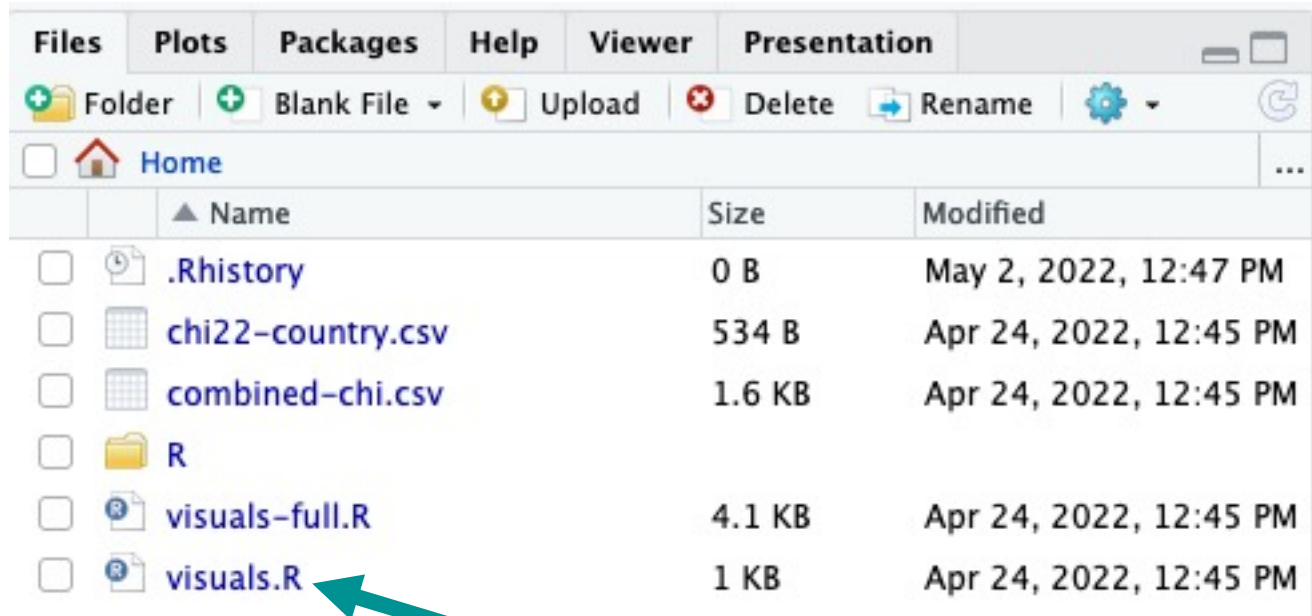
We're going to be visualizing some CHI data.

They're data about how publications by country in 2021 and 2022.

The 2021 and 2022 data are based on Kashyap Todi's releases.

Thanks Kashyap!

Loading CHI 2022 data



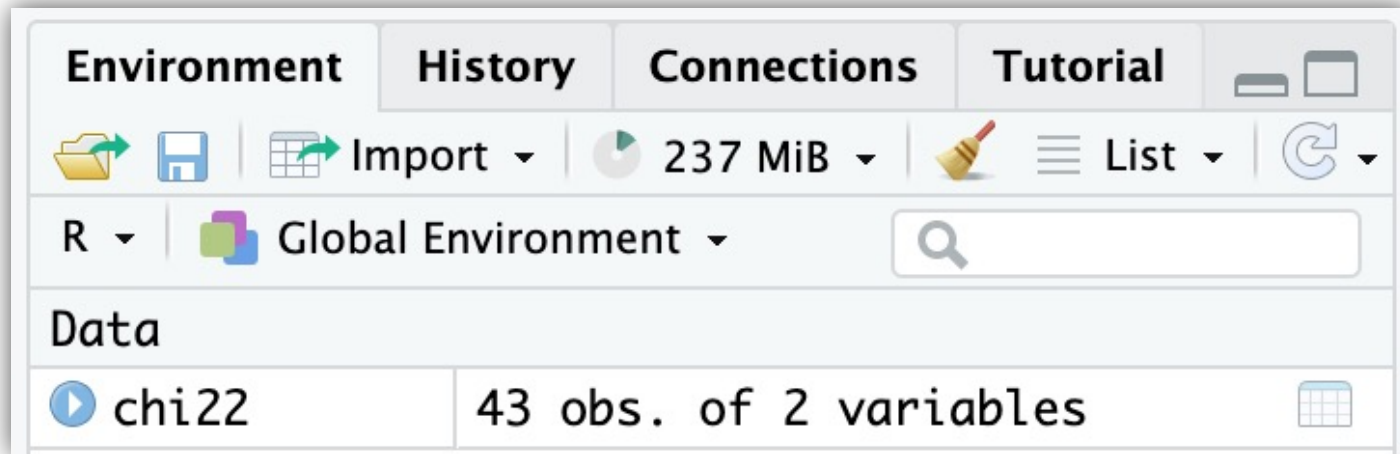
Click this!

CSV Import code

```
library("tidyverse")

##### Part1 A Simple Plot #####

#Import our 2022 country data
chi22 <- read.csv2("chi22-country.csv", sep=",")
```



The screenshot shows the RStudio Environment pane. At the top, there are tabs for 'Environment', 'History', 'Connections', and 'Tutorial'. Below the tabs is a toolbar with icons for file operations (folder, save, import), memory usage (237 MiB), a brush icon, a list icon, and a refresh icon. The main area shows 'R' and 'Global Environment' with a search bar. Below this, the 'Data' section is visible, showing a data object named 'chi22' with 43 observations of 2 variables.

Data	
chi22	43 obs. of 2 variables

CHI 2022 Country Data

n	Country	Pubs
1	United States	370
2	United Kingdom	97
3	Germany	69
4	Canda	63
5	China	42
6	Australia	35
7	Republic of Korea	33
8	France	29
9	Denmark	22
10	Japan	18

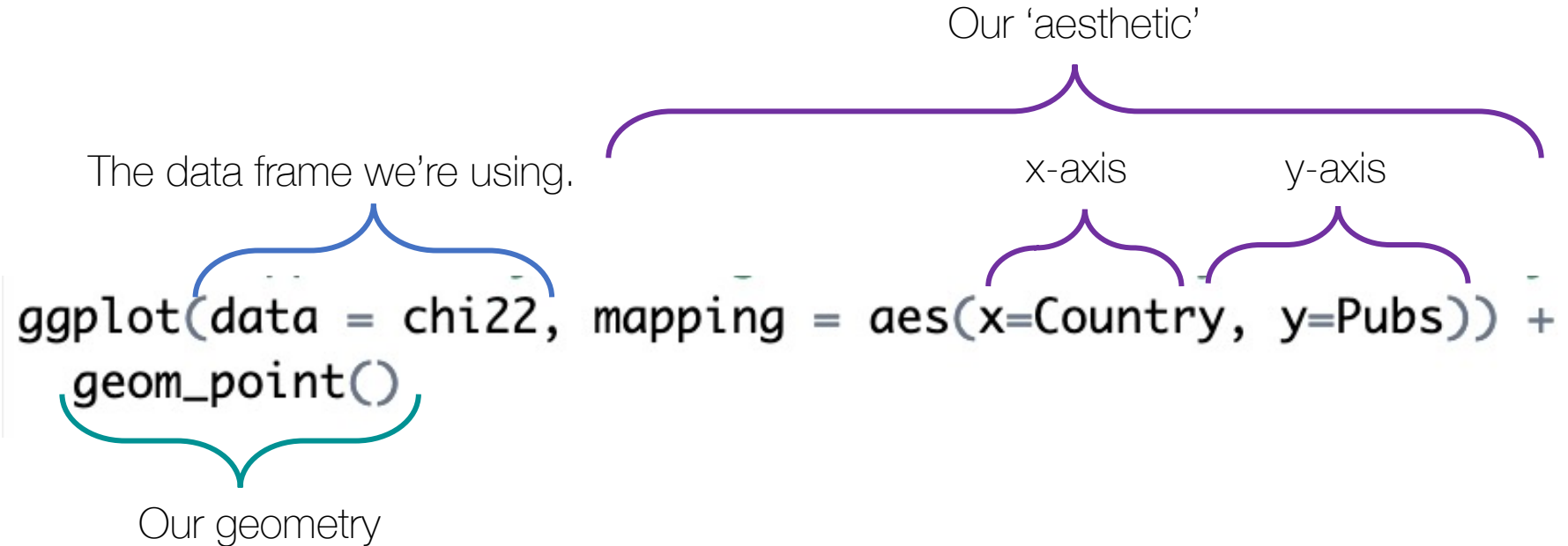
ggplot2

ggplot2 visualizations are created using the `ggplot()` function.

This function defines:

- The data we are using
- The aesthetics of the visualization we will create

The ggplot2 function



ggplot2 geometries



<code>geom_abline()</code>	<code>stat_sum()</code>	<code>geom_errorbar()</code>	<code>geom_segment()</code>
<code>geom_hline()</code>	<code>geom_density()</code>	<code>geom_linerange()</code>	<code>geom_curve()</code>
<code>geom_vline()</code>	<code>stat_density()</code>	<code>geom_pointrange()</code>	<code>geom_smooth()</code>
<code>geom_bar()</code>	<code>geom_density_2d()</code>	<code>geom_map()</code>	<code>stat_smooth()</code>
<code>geom_col()</code>	<code>stat_density_2d()</code>	<code>geom_path()</code>	<code>geom_spoke()</code>
<code>stat_count()</code>	<code>geom_dotplot()</code>	<code>geom_line()</code>	<code>geom_label()</code>
<code>geom_bin2d()</code>	<code>geom_errorbarh()</code>	<code>geom_step()</code>	<code>geom_text()</code>
<code>stat_bin_2d()</code>	<code>geom_hex()</code>	<code>geom_point()</code>	<code>geom_raster()</code>
<code>geom_blank()</code>	<code>stat_bin_hex()</code>	<code>geom_polygon()</code>	<code>geom_rect()</code> <code>geom_tile()</code>
<code>geom_boxplot()</code>	<code>geom_freqpoly()</code>	<code>geom_qq_line()</code>	<code>geom_violin()</code>
<code>stat_boxplot()</code>	<code>geom_histogram()</code>	<code>stat_qq_line()</code>	<code>stat_ydensity()</code>
<code>geom_contour()</code>	<code>stat_bin()</code>	<code>geom_qq()</code> <code>stat_qq()</code>	<code>stat_sf()</code> <code>geom_sf()</code>
<code>stat_contour()</code>	<code>geom_jitter()</code>	<code>geom_quantile()</code>	<code>geom_sf_label()</code>
<code>geom_count()</code>	<code>geom_crossbar()</code>	<code>stat_quantile()</code>	<code>geom_sf_text()</code>
<code>geom_ribbon()</code>	<code>geom_area()</code>	<code>geom_rug()</code>	<code>coord_sf()</code>

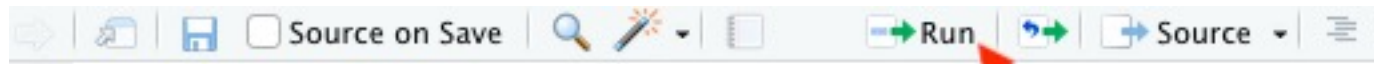
The ggplot2 function

Now select the code and run it:

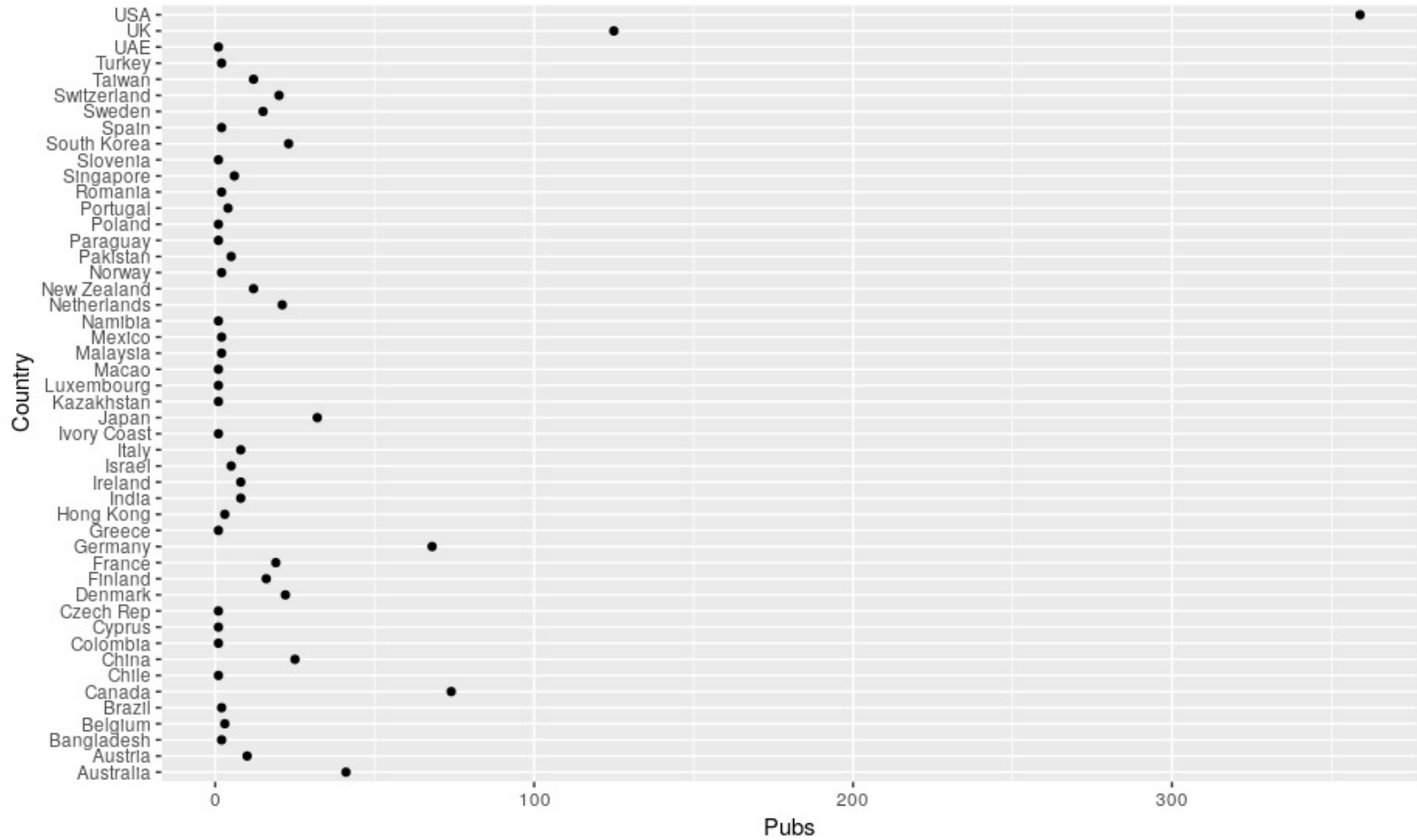
Swap x and y



```
ggplot(data = chi22, mapping = aes(x=Country, y=Pubs)) +  
  geom_point()
```



Click this!



Let's try and build this out a bit

You should now have a basic plot. It's far from perfect though

Let's pick up the webpage at **5.3 Getting some order to proceedings**

You're going to:

- Get the data into a sensible order
- Have a play around with different geometries

Stop when you get to **6. Controlling the appearance of plots**

Controlling the look of plots

There are three ways of controlling the look of plots:

- Through aesthetics

- Through geometries

- Through themes



We'll look at these first

Appearance in geometries

```
ggplot(data = chi22 mapping = aes(x=reorder(Country, size), y=chi22)) +  
  geom_point(colour="Orange", size=3) +  
  coord_flip()
```

Color of points

Size of points

Appearance in geometries

On the website, work from

6.1 Appearance as part of geometries

Until you reach the end of:

6.2 Controlling appearance with themes

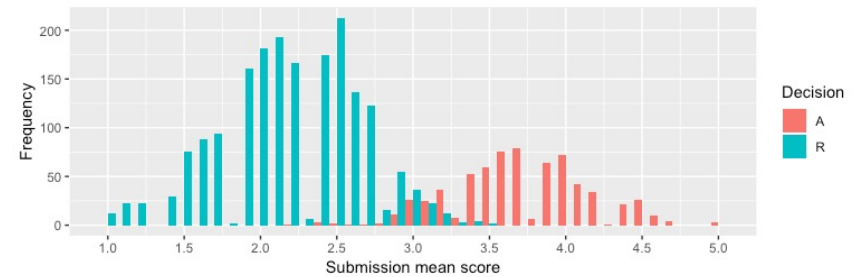
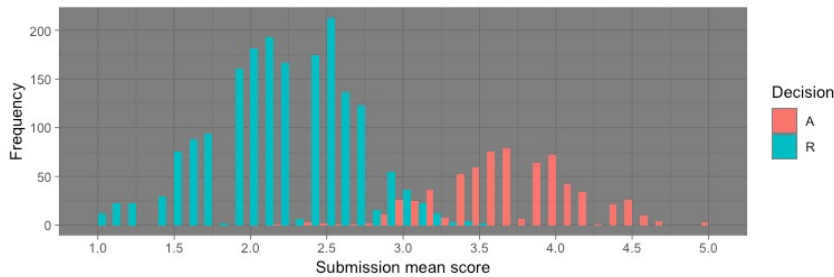
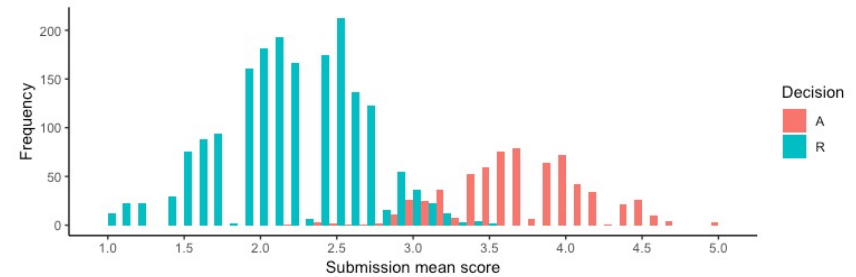
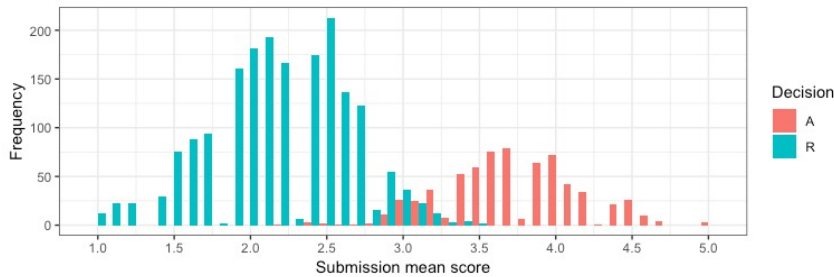
Controlling appearance with themes

Themes control the **overall** appearance of your plot.

- Gridlines
- Background
- Axes
- Ticks
- Labels

Controlling appearance with themes

There are built-in themes, like `theme_classic()`, `theme_dark()`, `theme_bw()` and `theme_gray()`.



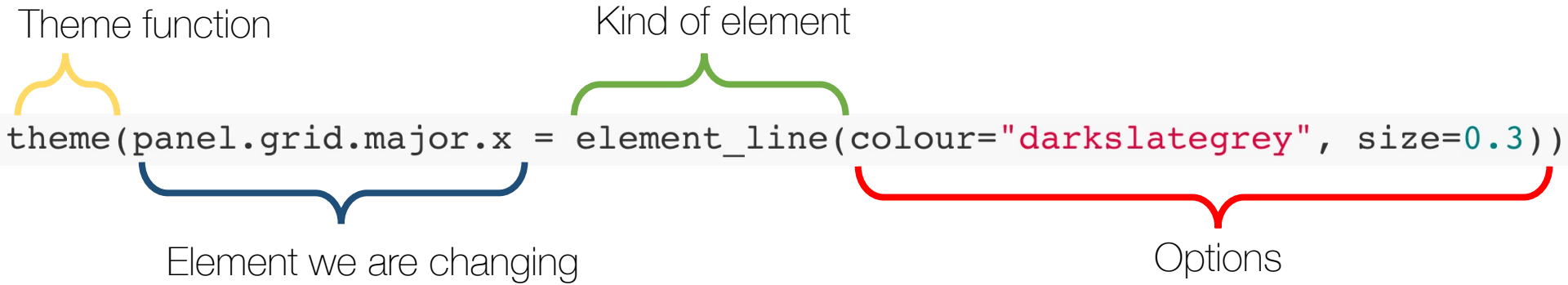
Controlling appearance with themes

Themes are made up of a number of components in a hierarchical form. (e.g., `legend.box.background` or `plot.caption`).

They are drawn using elements:

- Lines are `element_line()`
- Rectangles are `element_rect()`
- etc

Controlling appearance with themes



Controlling appearance with themes

On the website, work until you reach:

7. A plot drawing on two datasets

Plotting multiple variables

So far we have looked at a very simple plot.

- It has a categorical variable.
- And a continuous variable.

What if we want to look at many year's worth of data?


Enter the **CHI**mbined dataset with 2021 and 2022 data!

CHI 2021 and 2022 Country Data

n	Country	Year	Pubs
1	Australia	2021	50
2	Australia	2022	35
3	Austria	2021	11
4	Austria	2022	5
5	Bangladesh	2021	4
6	Bangladesh	2022	2
7	Barbados	2021	1
8	Belgium	2021	5
9	Belgium	2022	6
10	Brazil	2021	3

Plotting multiple variables

There are three ways of controlling the look of plots:

- Through aesthetics 
- Through geometries
- Through themes

Plotting multiple variables

aes function

New 'colour' aes addition

```
aes(x=reorder(Country, Pubs), y=Pubs, colour=factor(Year))
```

Stuff you recognize

Plotting multiple variables

On the website, work from

7. A plot drawing on multiple variables

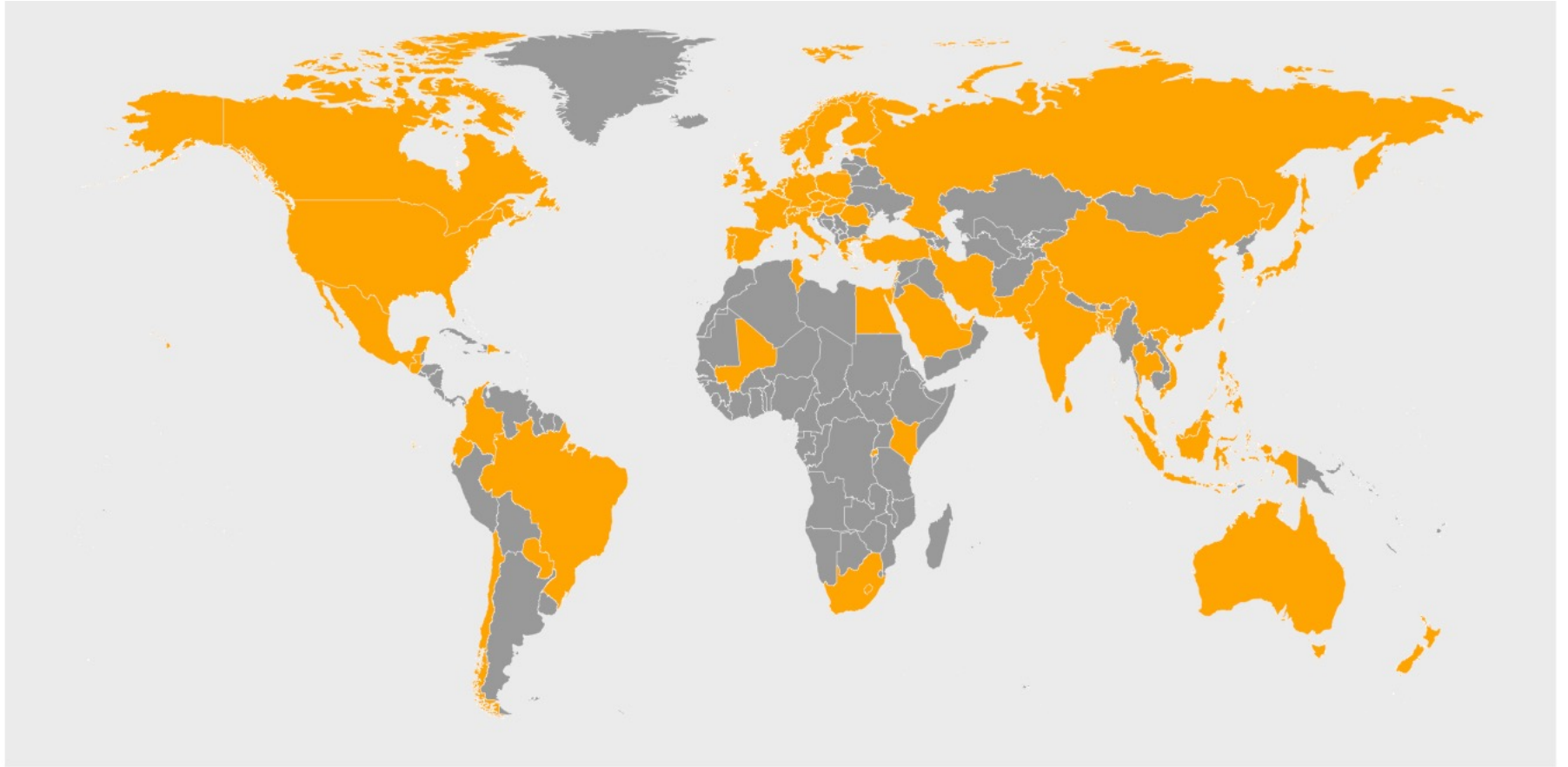
Until you reach:

8. Mapping locations of CHI authors

Mapping locations

```
#Import mapping libraries  
library(maps)
```

Mapping locations



Mapping locations

```
world <- map_data("world")  
world$fac <- world$region %in% chi22 $Country
```

Mapping locations

We're going to use the `geom_poly()*` geometry for plotting our graph.

Map of the world Longitude and latitude

```
ggplot(data = world, aes(x=long, y = lat, group = group)) +  
  geom_polygon(colour="white")
```

Plot the polygons Colour borders in white Keep countries together

* There is a `geom_map()`, which I discuss on the website

Mapping locations

We're going to use the `geom_poly()*` geometry for plotting our graph.

Map of the world Longitude and latitude

```
ggplot(data = world, aes(x=long, y = lat, group = group)) +  
  geom_polygon(colour="white")
```

Plot the polygons Colour borders in white Keep countries together

* There is a `geom_map()`, which I discuss on the website

Plotting multiple variables

On the website, work from

8. Mapping locations of CHI authors

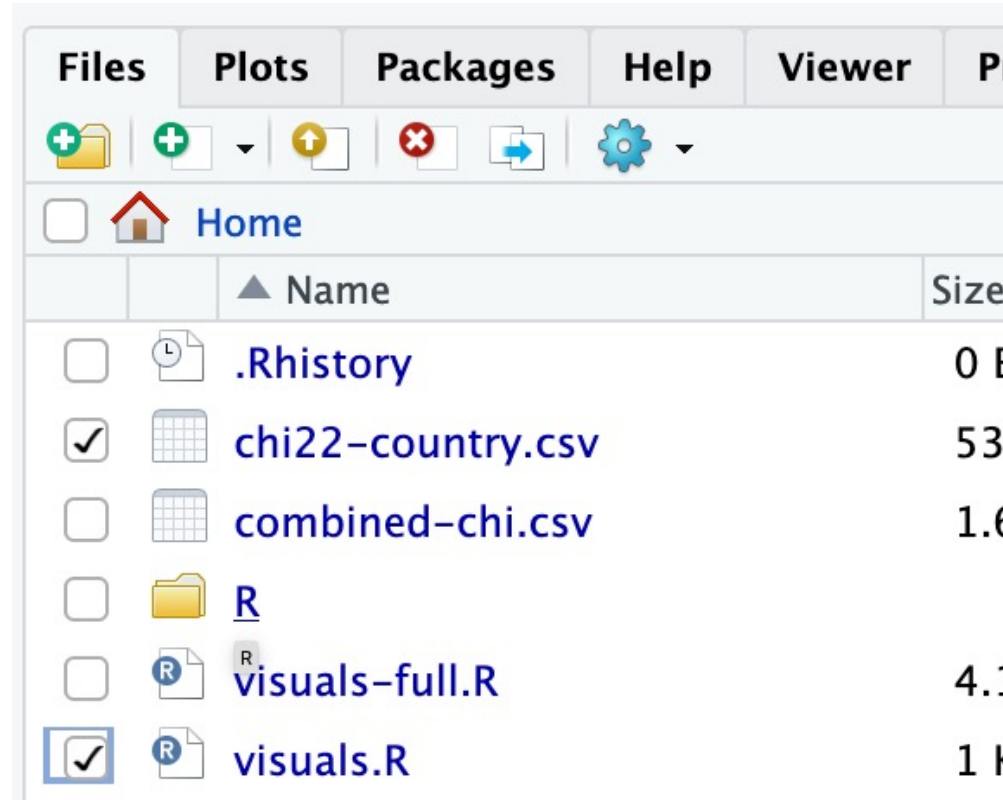
Until you reach:

8. We're done... for now!

Getting your code

Download your files before you lose them!

Select the files you want to keep →



Getting your code

Download your files before you lose them!

**Then click
More →
Export...**

**To get a zip
file.**

