

Coding data science

March 2019

Hadley Wickham

@hadleywickham

Chief Scientist, RStudio



What is data analysis?

Data analysis is the process
by which data becomes
understanding, knowledge
and insight

Data analysis is the process
by which data becomes
understanding, knowledge
and insight

Import



Tidy



Consistent way of
storing data

Import



Tidy

Consistent way of
storing data



Understand

Import



Tidy

Consistent way of
storing data



Transform

Create new variables & new summaries

Visualise

Surprises, but doesn't scale



Model

Scales, but doesn't (fundamentally) surprise



Import



Tidy

Consistent way of
storing data



Transform

Create new variables & new summaries

Visualise

Surprises, but doesn't scale

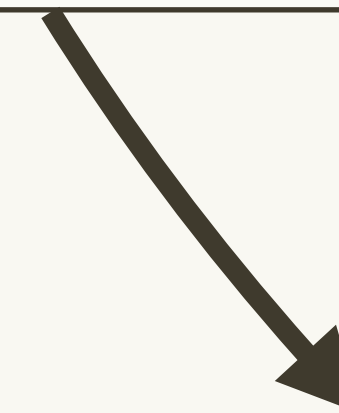


Model

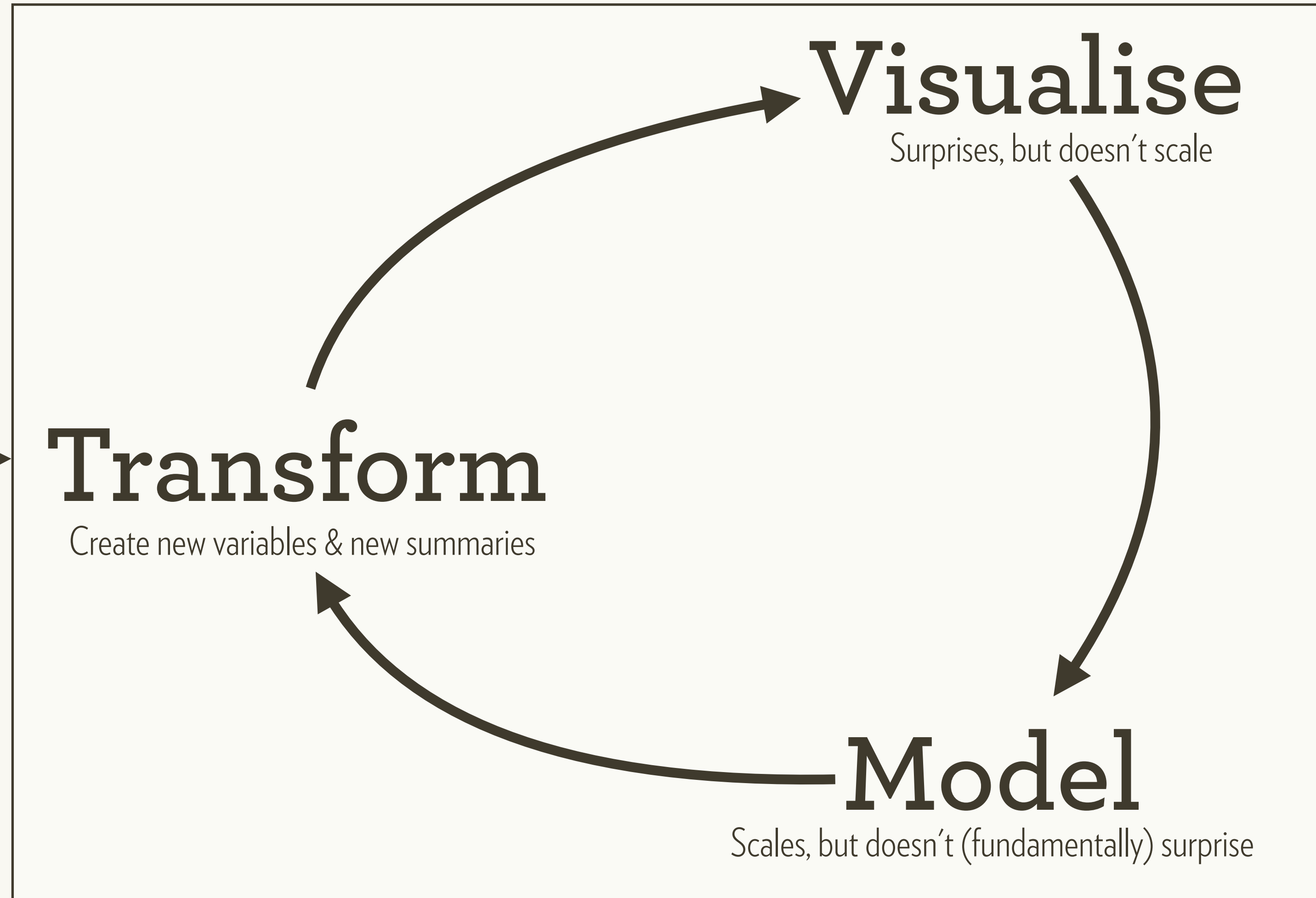
Scales, but doesn't (fundamentally) surprise



Communicate



Automate



What is data science?

Data science =
data analysis +
programming

Import



Tidy

Consistent way of
storing data

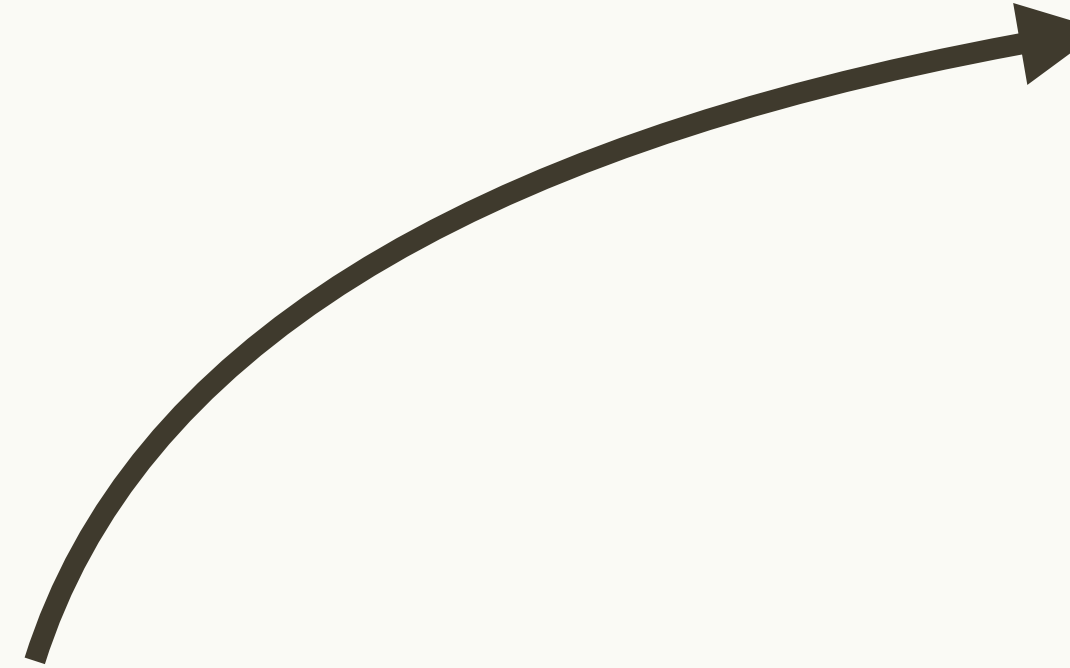


Transform

Create new variables & new summaries

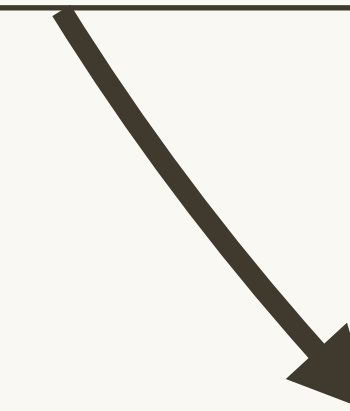
Visualise

Surprises, but doesn't scale



Model

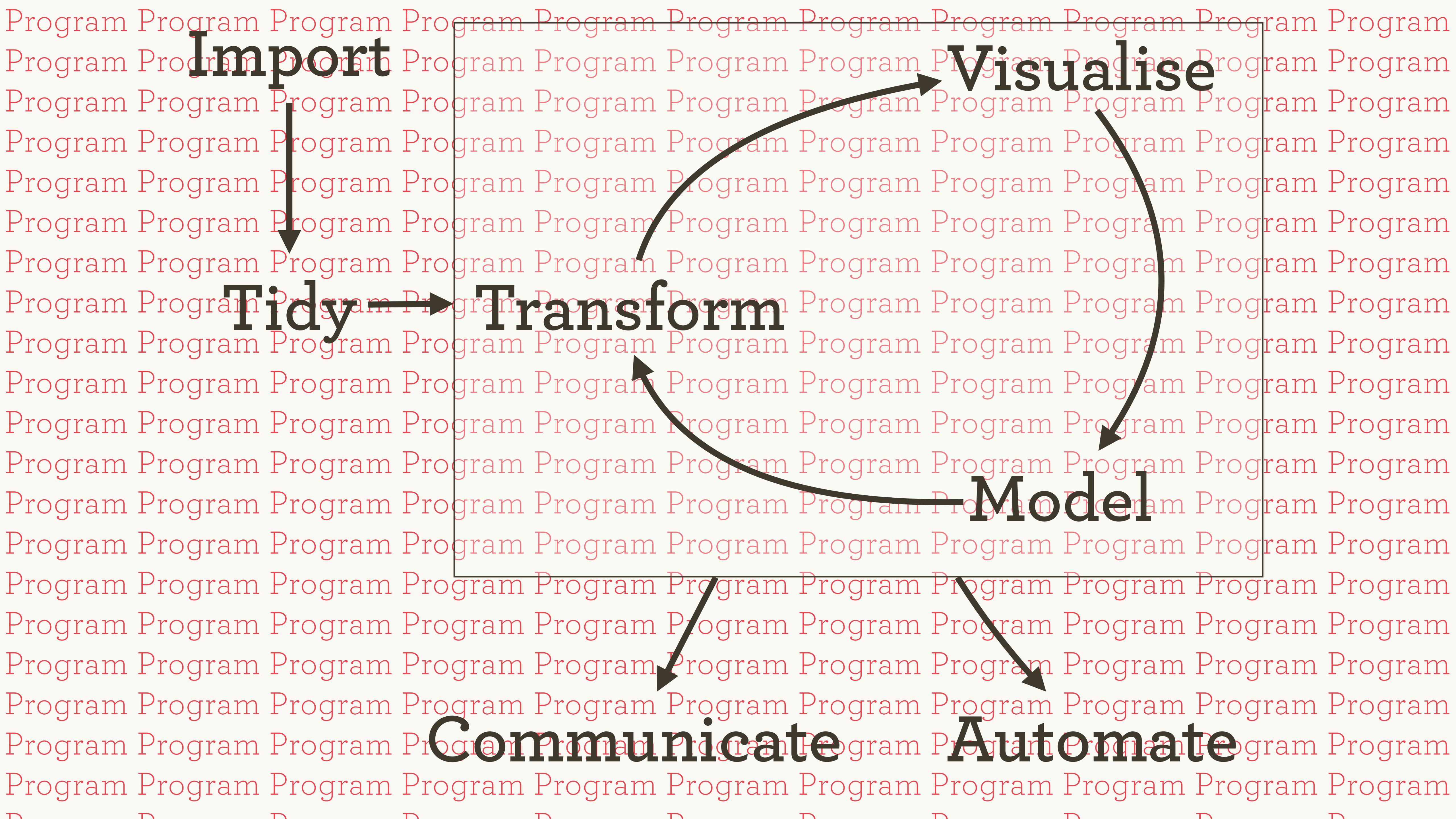
Scales, but doesn't (fundamentally) surprise



Program

Communicate

Automate





Import

readr
readxl
haven
xml2

Tidy

tibble
tidyr

Transform

dplyr
forcats
hms

lubridate
stringr

Visualise

ggplot2

Model

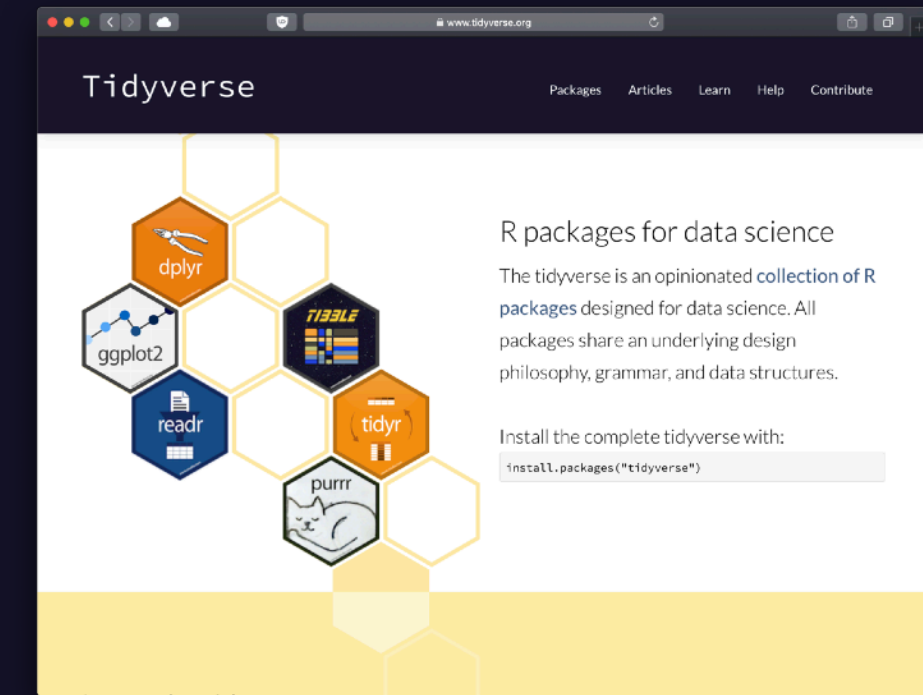
recipes
parsnip

purrr
magrittr

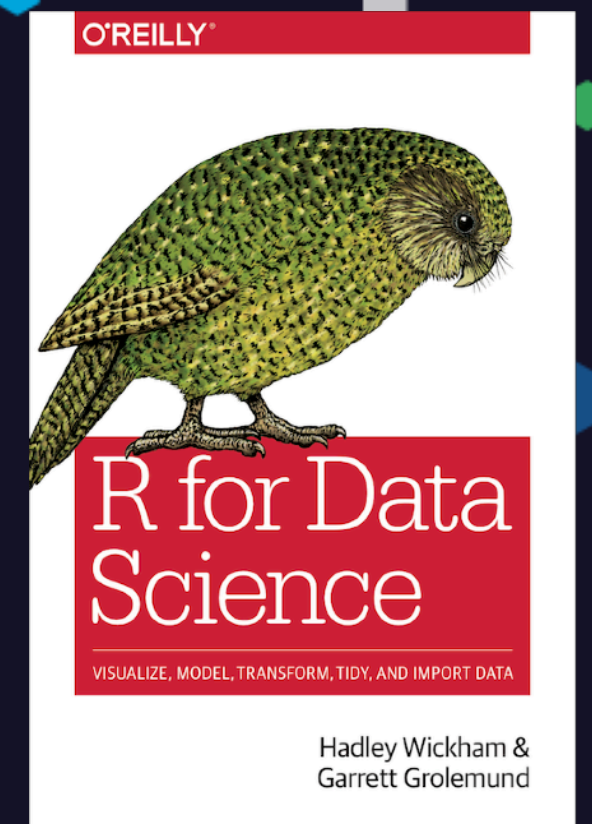
Program

shiny
rmarkdown

Communicate



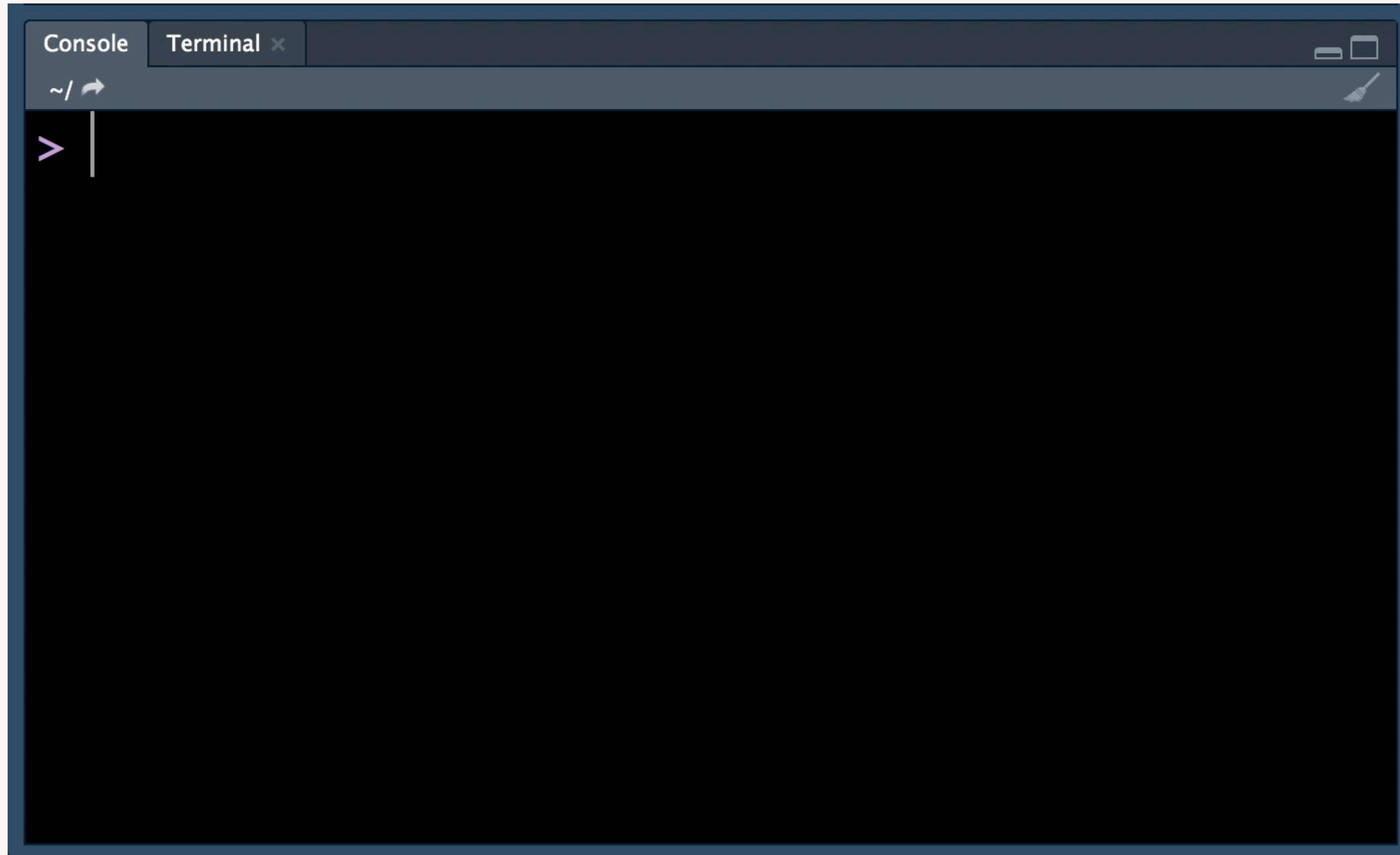
tidyverse.org



r4ds.had.co.nz

Why code?

The disadvantages of code are obvious

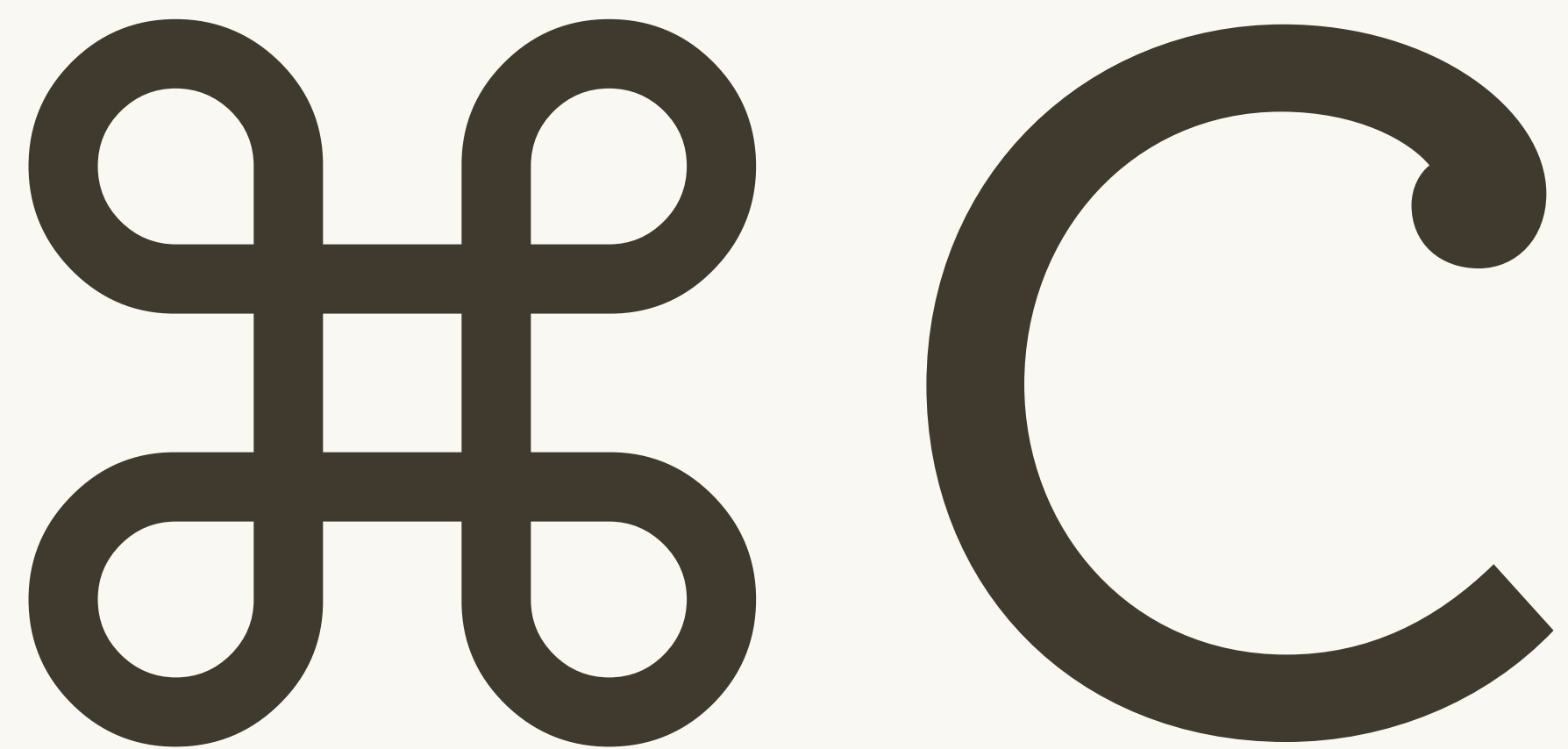


Why code?

1. Code is text

2. Code is read-able

3. Code is reproducible



Copy



Paste

QuestionsDeveloper JobsTagsUsers

Search...

4+2189

Top Questions

interesting391featuredhotweekmonth

0

0

2

votesanswersviews

Where to scroll to in a react app when route changes for screen reader

reactjsreact-routeraccessibilityreact-router-dom

asked 55 secs ago dagda19,694

0

0

2

votesanswersviews

what if i schedule tasks for celery to perform every minute and it is not able to complete it in time?

celeryscheduled-taskscelery-taskcelerybeat

asked 1 min ago ravi1

1

0

5

voteanswersviews

Gerrit: Is there a way to push directly into master?

gerrit

modified 2 mins ago leeyuiwah1,909

0

0

3

votesanswersviews

kubectl run-ing a tarball'd image

dockerkuberneteskubectl

asked 2 mins ago adelbertc4,470

0

0

7

votesanswersviews

AspNet Identity RequireUniqueEmail = false throws exception on CreateAsync

asp.net-identityuniqueowin

modified 2 mins ago Dmitry Duka1

0

0

3

votesanswersviews

Square: Call to undefined function charge

square-connect

asked 3 mins ago John4,104

0

0

4

votesanswersviews

Python - Social Studio API - How to unpack JSON into pandas data frame

pythonjsonpandas

asked 3 mins ago Ulises Sotomayor36

2

1

15

votesanswerviews

Issues running airflow scheduler as a daemon process

pythonamazon-ec2ubuntu-16.04airflowapache-airflow

answered 3 mins ago Tagar2,593

1

1

20

voteanswerviews

Flutter animation how to fade in/out gradually

flutter

answered 3 mins ago Collin Jackson6,957

Why code?

1. Code is text

2. Code is read-able

3. Code is reproducible

What have you done?

The image shows a Microsoft Excel spreadsheet with a 'Sort Warning' dialog box open. The spreadsheet has columns A through L and rows 1 through 24. The data in columns A through E is as follows:

	A	B	C	D	E
1	a	b	c	d	e
2	0.78	0.91	0.56	0.21	0.74
3	0.30	0.53	0.91	0.50	0.74
4	0.25	0.57	0.56	0.05	0.86
5	0.81	0.19	0.99	0.63	0.42
6	0.90	0.37	0.55	0.55	0.49
7	0.14	0.69	0.84	0.41	0.65
8	0.68	0.08	0.28	0.08	0.19
9	0.77	0.87	0.66	0.92	0.71
10	0.71	0.00	0.90	0.06	0.95
11	0.16	0.53	0.09	0.34	0.17
12	0.17	0.86	0.86	0.48	0.94
13	0.12	0.32	0.15	0.43	0.21
14	0.64	0.78	0.85	0.06	0.69
15	0.78	0.23	0.47	0.31	0.66
16	0.87	0.31	0.28	0.16	0.74
17	0.68	0.42	0.66	0.23	0.13
18	0.02	0.18	0.03	0.41	0.64
19	0.44	0.65	0.87	0.09	0.09
20	0.67	0.54	0.28	0.20	0.99
21	0.74	0.03	0.58	0.61	0.68
22					
23					
24					

The 'Sort Warning' dialog box is open, asking 'What do you want to do?'. The options are:

- ☒ Expand the selection
- ☐ Continue with the current selection

The 'Sort' button is highlighted.

Why code?

1. Code is text

2. Code is read-able

3. Code is reproducible

```

# install.packages("devtools")
devtools::install_github("jennybc/frogs")

## Getting to know the frogs

At this point, all we know is that each row is one frog-jump. Frog ids coming ...

{r}
library(frogs)
library(tidyverse)

frogs
glimpse(frogs)

An early figure. Do frogs need to warm up? Do they fatigue? Yes and yes.

{r frog-fatigue, echo = FALSE}
frogs2 <- frogs %>%
  filter(jump_n < 7) %>%
  mutate(
    jump_n = as.factor(as.integer(jump_n))
  )
ggplot(frogs2, aes(x = distance, color = jump_n)) +
  geom_density()

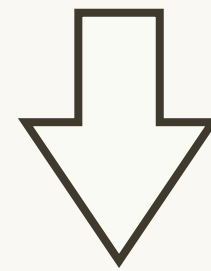
Do professional frog jumping teams get better results? YES.

{r frog-type, echo = FALSE}
ggplot(frogs, aes(x = distance, color = frog_type)) +
  geom_density()

```

.Rmd

Prose and code



The screenshot shows a web browser window with the GitHub repository for the 'frogs' R package. The page title is 'Getting to know the frogs'. The content is as follows:

At this point, all we know is that each row is one frog-jump. Frog ids coming ...

```
library(frogs)
library(tidyverse)

frogs
glance(frogs)
```

An early figure. Do frogs need to warm up? Do they fatigue? Yes and yes.

```
frogs2 <- frogs %>%
  filter(jump_n < 7) %>%
  mutate(
    jump_n = as.factor(as.integer(jump_n))
  )
ggplot(frogs2, aes(x = distance, color = jump_n)) +
  geom_density()
```

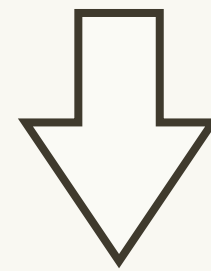
Do professional frog jumping teams get better results? YES.

```
ggplot(frogs, aes(x = distance, color = frog_type)) +
  geom_density()
```

At the bottom of the page, there is a footer with the text: © 2018 GEHUB, INC. Terms Privacy Security Status Help. On the right side of the footer, there are links: Contact GitHub API Training Shop Blog About.

.md

Prose and results



```
library(frogs)
library(tidyverse)

#> # ggplot2 2.2.1           Date: 2017-05-24
#> # tibble 1.2.1            R: 3.2.2
#> # tidyr 0.6.2.9000        OS: OS X El Capitan 10.11.6
#> # readr 1.1.0             GUI: X11
#> # purrr 0.2.2.9000        locales en_CA.UTF-8
#> # dplyr 0.6.0             TZ: America/Vancouver
#> # stringr 1.2.0
#> # forcats 0.2.0
#> # Conflicts —————
#> # filter(), from dplyr, masks stats::filter()
#> # lag(), from dplyr, masks stats::lag()
```

```
frogs
#> # A tibble: 3,272 x 15
#>   row distance distance_3_jump_n frog_type distance_3_off
#>   <int> <dbl> <dbl> <dbl> <int> <chr> <dbl>
#> 1     1    165.900    0.58333      0    1 pro -1
#> 2     2    177.400    0.71667      2    2 pro -1
#> 3     3     0.000    0.00000      0    3 pro -1
#> 4     4    27.158    0.43333      0    1 pro -1
#> 5     5     0.000    0.00000      2    2 pro -1
#> 6     6     0.000    0.00000      0    3 pro -1
#> 7     7    48.914    0.48889      0    1 pro -1
#> 8     8     0.000    0.00000      2    2 pro -1
#> 9     9     0.000    0.00000      3    3 pro -1
#> 10    10 35.853    0.43333      0    1 pro -1
#> # ... with 3,262 more rows, and 0 more variables: distance_rel <dbl>,
#> # day <dbl>, angle_01 <dbl>, angle_10 <dbl>, angle_00 <dbl>,
#> # velocity_01<dbl>, velocity_10<dbl>, velocity_00 <dbl>
glimpse(frogs)
#> Observations: 3,272
#> Variables: 15
```

.html

Human shareable


```

# install.packages("devtools")
devtools::install_github("jennybc/frogs")

## Getting to know the frogs

At this point, all we know is that each row is one frog-jump. Frog ids coming ...

```{r}
library(frogs)
library(tidyverse)

frogs
glimpse(frogs)
```

An early figure. Do frogs need to warm up? Do they fatigue? Yes and yes.

```{r frog-fatigue, echo = FALSE}
frogs2 <- frogs %>%
 filter(jump_n < 7) %>%
 mutate(
 jump_n = as.factor(as.integer(jump_n))
)
ggplot(frogs2, aes(x = distance, color = jump_n)) +
 geom_density()
```

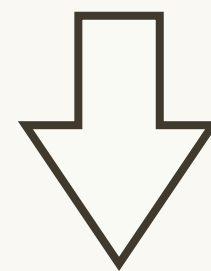
Do professional frog jumping teams get better results? YES.

```{r frog-type, echo = FALSE}
ggplot(frogs, aes(x = distance, color = frog_type)) +
 geom_density()
```

```

.Rmd

Prose and code



Getting to know the frogs

At this point, all we know is that each row is one frog-jump. Frog ids coming ...

```
library(frogs)
library(tidyverse)

frogs
glance(frogs)
```

An early figure. Do frogs need to warm up? Do they fatigue? Yes and yes.

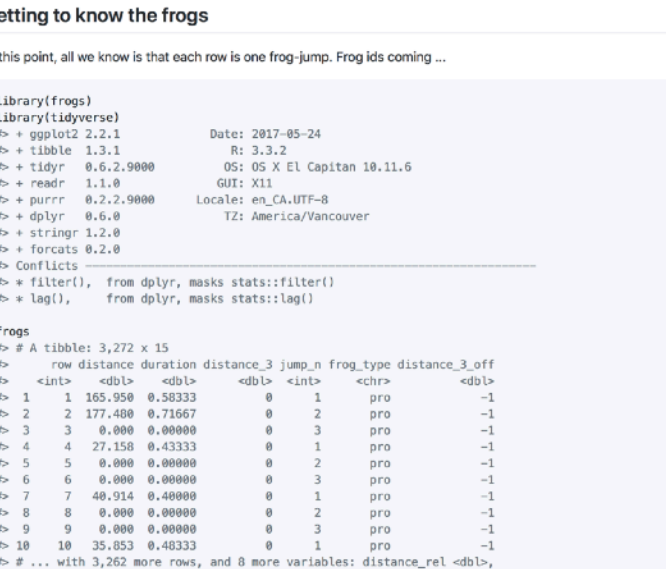
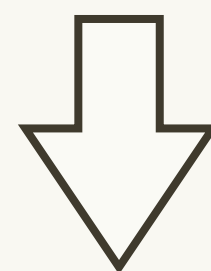
```
frogs2 <- frogs %>%
  filter(jump_n < 7) %>%
  mutate(
    jump_n = as.factor(as.integer(jump_n))
  )
ggplot(frogs2, aes(x = distance, color = jump_n)) +
  geom_density()
```

Do professional frog jumping teams get better results? YES.

```
ggplot(frogs, aes(x = distance, color = frog_type)) +
  geom_density()
```

.md

Prose and results



```

library(frogs)
library(tidyverse)
#> ggplot2 2.2.1          Date: 2017-05-24
#> tibble 1.1.1           R: 3.3.2
#> tidyr 0.6.2.9000      OS: OS X El Capitan 10.11.6
#> readr 1.1.0           GUI: X11
#> purrr 0.2.2.9000      Locale: en_CA.UTF-8
#> dplyr 0.6.0           TZ: America/Vancouver
#> string 1.2.0
#> forcats 0.2.0

#> Conflicts —————
#> =========
#> w = filter(), from dplyr, masks stats::filter()
#> = lag(), from dplyr, masks stats::lag()

frogs
#> # A tibble: 3,272 x 15
#>   row distance distance_3_jump_n frog_type distance_3_off
#>   <int> <dbl> <dbl> <dbl> <int> <chr> <dbl>
#> 1 1 165.950 0.50333 0 1 1 pro -1
#> 2 2 177.400 0.71607 0 2 pro -1
#> 3 3 0.000 0.00000 0 3 pro -1
#> 4 4 27.158 0.43333 0 1 pro -1
#> 5 5 0.000 0.00000 0 2 pro -1
#> 6 6 0.000 0.00000 0 3 pro -1
#> 7 7 48.914 0.40000 0 1 pro -1
#> 8 8 0.000 0.00000 0 2 pro -1
#> 9 9 0.000 0.00000 0 3 pro -1
#> 10 10 35.853 0.40333 0 1 pro -1
#> # ... with 3,262 more rows, and 0 more variables: distance_rel_dbl,
#>   day_dbl, angle_at_dbl, angle_10_dbl, angle_80_dbl,
#>   velocity_0_dbl, velocity_10_dbl, velocity_80_dbl
glimpse(frogs)
#> Observations: 3,272
#> Variables: 15

```

.html

.doc

.tex

.ppt

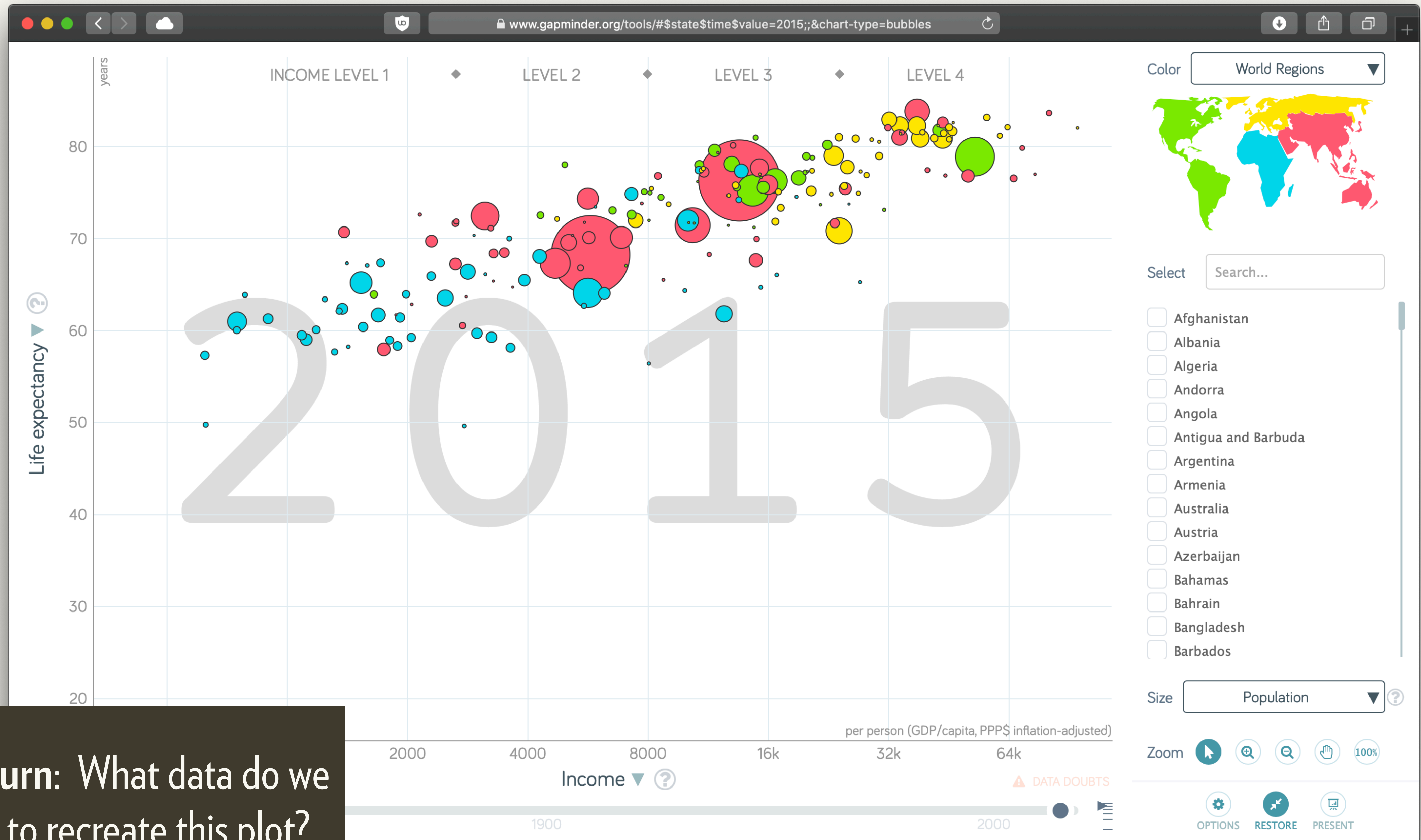
.pdf

● ● ●

What about non-
programmers?

You don't need to
be a programmer
to code!

Your turn: What data do we need to recreate this plot?



Underlying data

```
# A tibble: 193 x 6
```

| | country | four_regions | year | income | life_exp | pop |
|----|---------------------|---------------------|-------------|---------------|-----------------|------------|
| | <chr> | <chr> | <int> | <int> | <dbl> | <int> |
| 1 | Afghanistan | asia | 2015 | 1750 | 57.9 | 33700000 |
| 2 | Albania | europa | 2015 | 11000 | 77.6 | 2920000 |
| 3 | Algeria | africa | 2015 | 13700 | 77.3 | 39900000 |
| 4 | Andorra | europa | 2015 | 46600 | 82.5 | 78000 |
| 5 | Angola | africa | 2015 | 6230 | 64 | 27900000 |
| 6 | Antigua and Barbuda | americas | 2015 | 20100 | 77.2 | 99900 |
| 7 | Argentina | americas | 2015 | 19100 | 76.5 | 43400000 |
| 8 | Armenia | europa | 2015 | 8180 | 75.4 | 2920000 |
| 9 | Australia | asia | 2015 | 43800 | 82.6 | 23800000 |
| 10 | Austria | europa | 2015 | 44100 | 81.4 | 8680000 |

```
# ... with 183 more rows
```

Phonics are important!

```
gapminder %>%
```

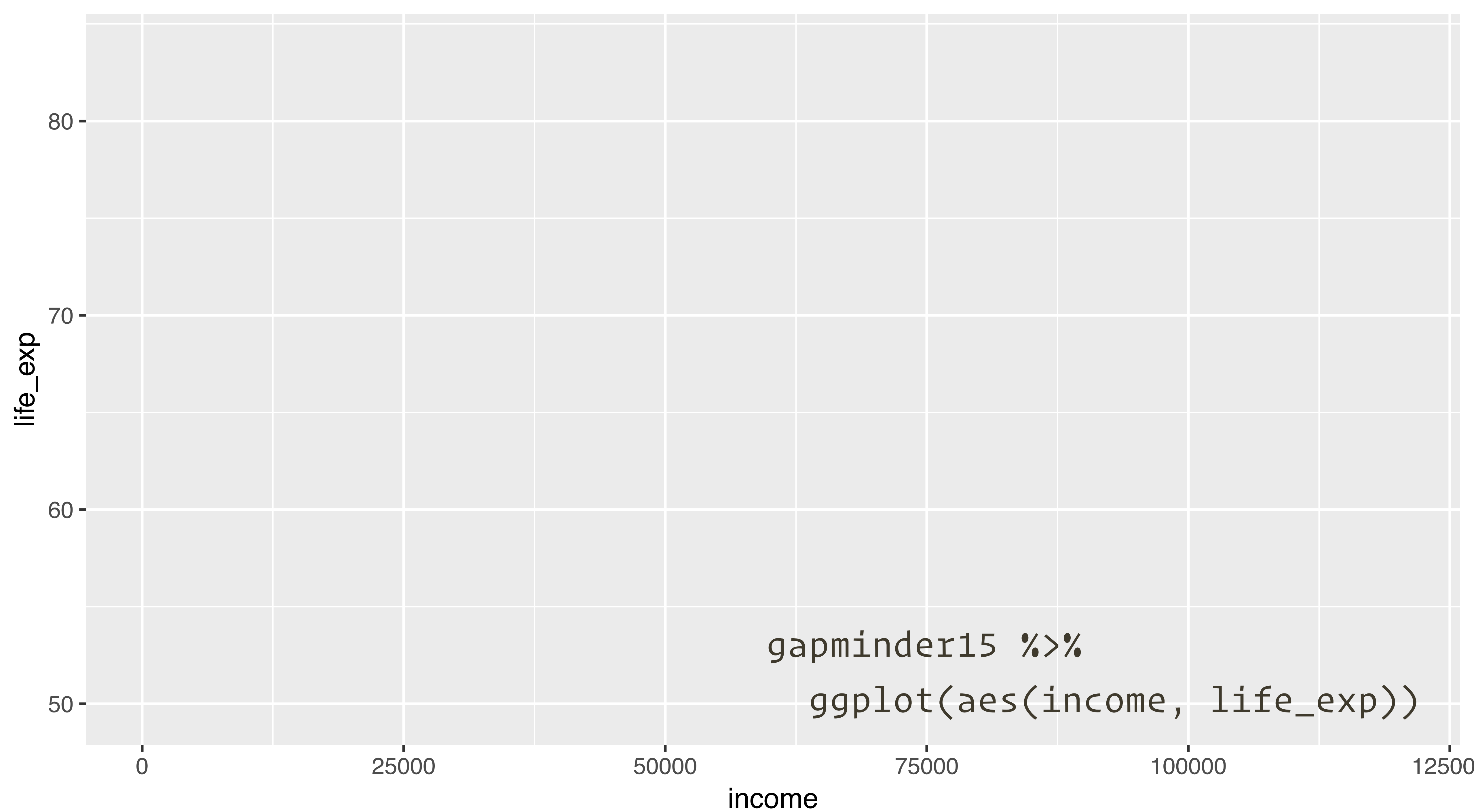
Take the gapminder data, then

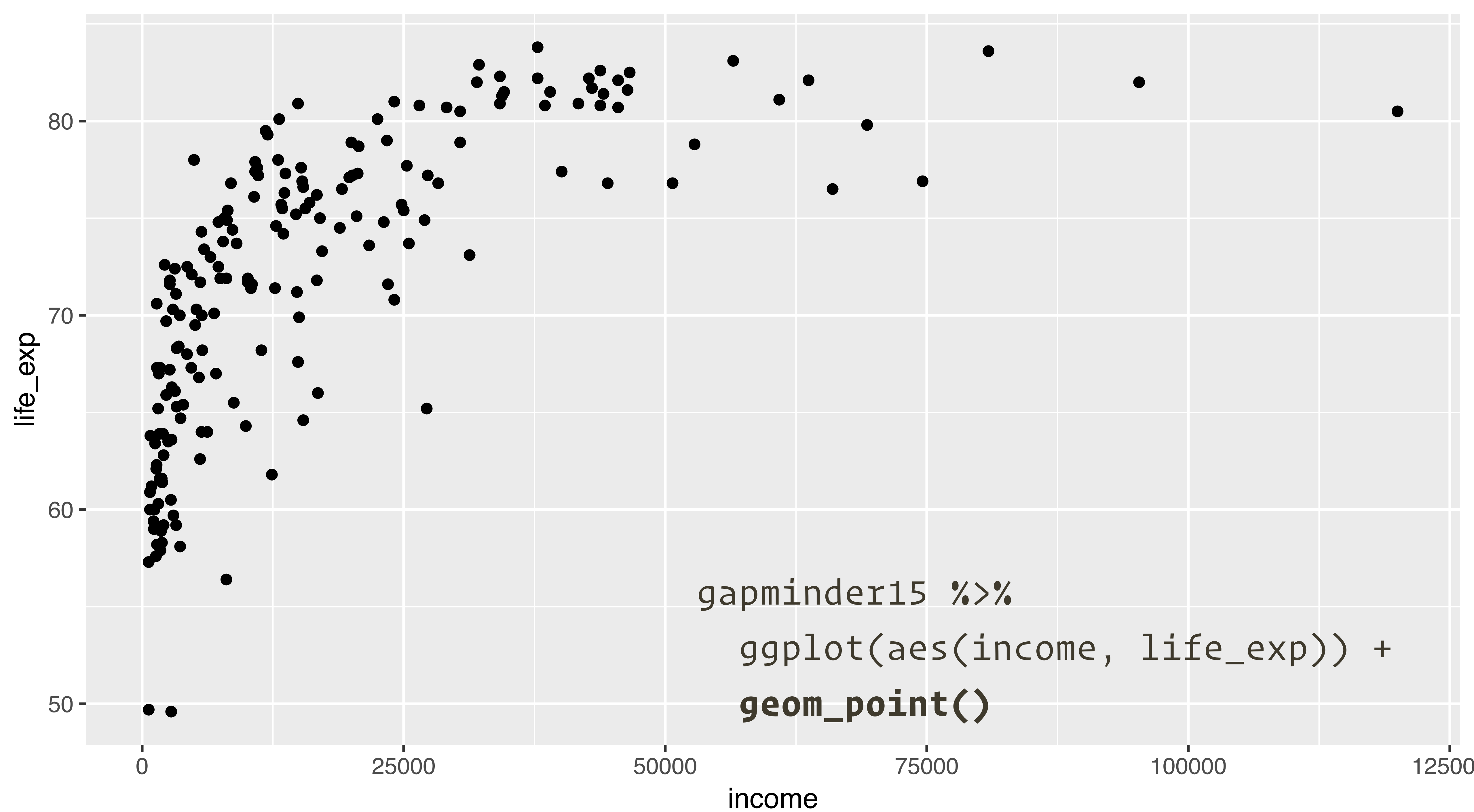
```
filter(year == 2015) ->
```

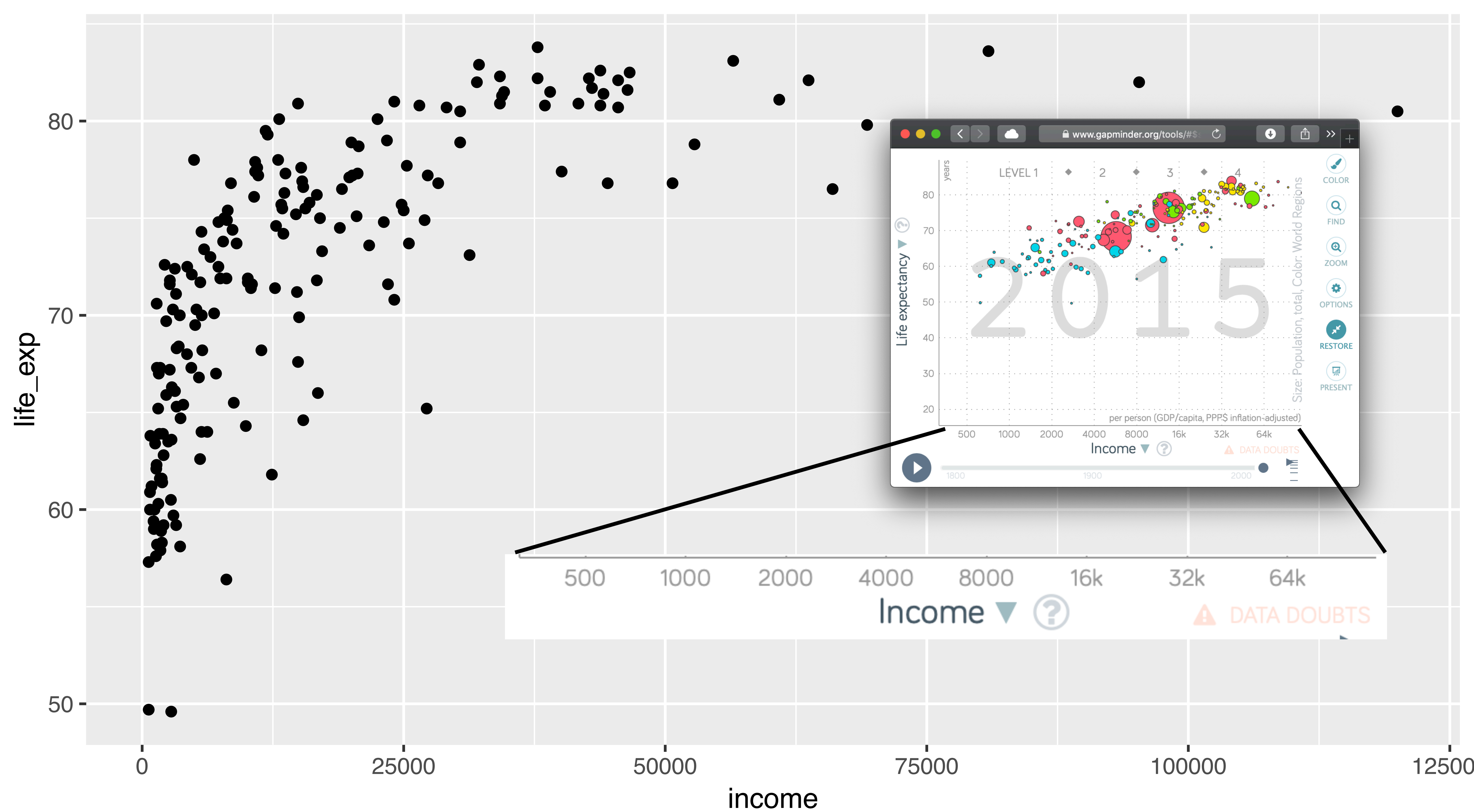
filter rows where year equals 2015, creating

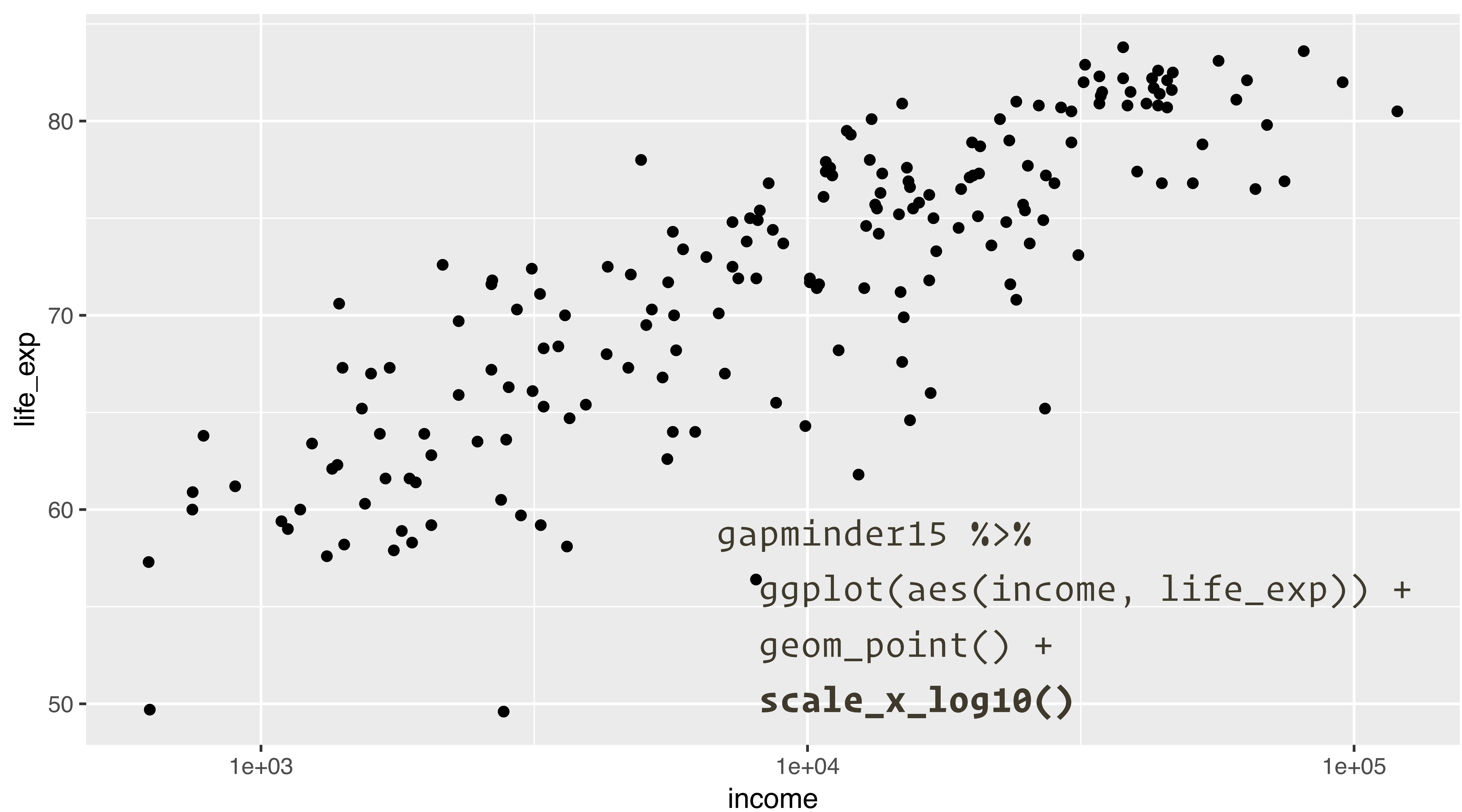
```
gapminder15
```

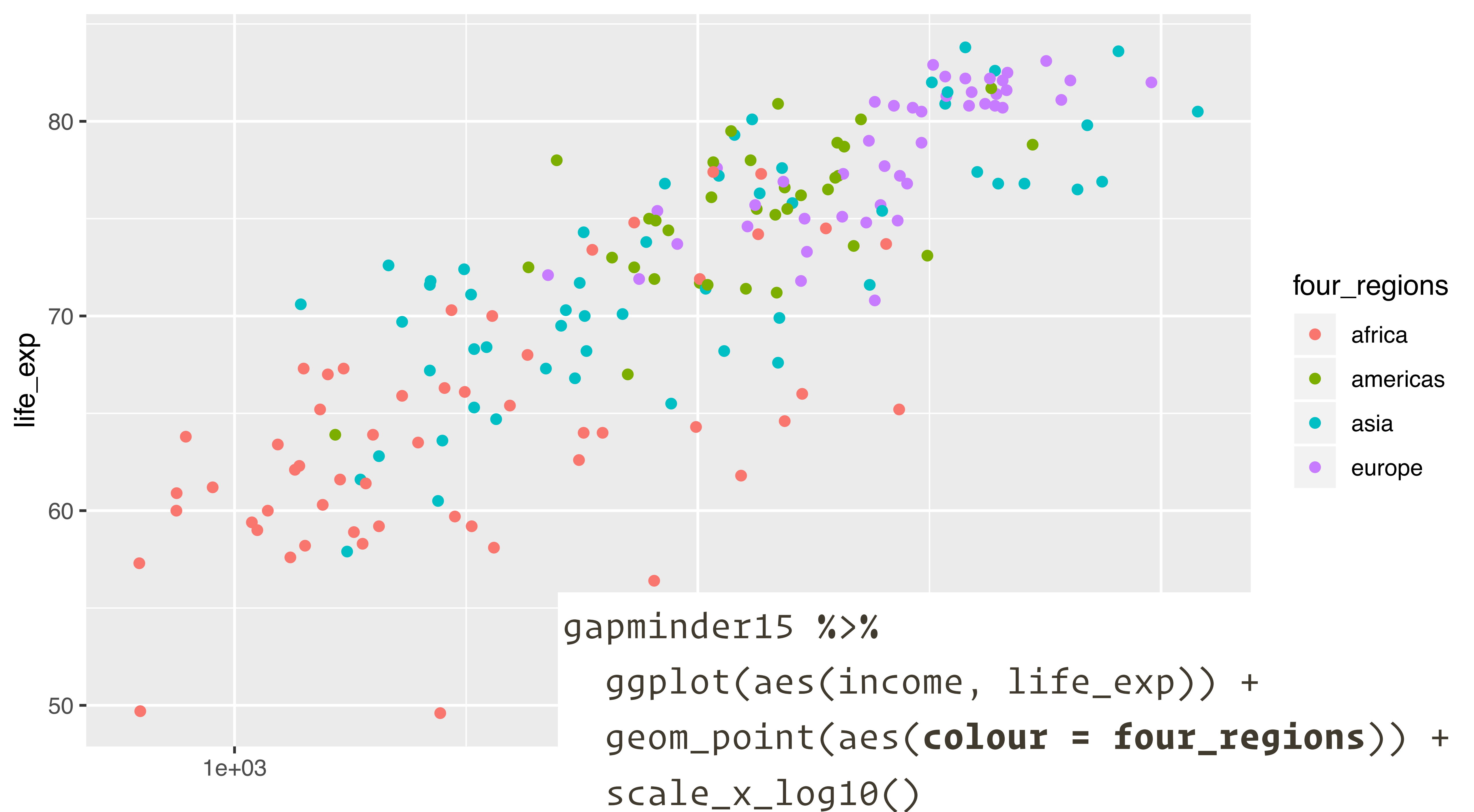
gapminder15 variable

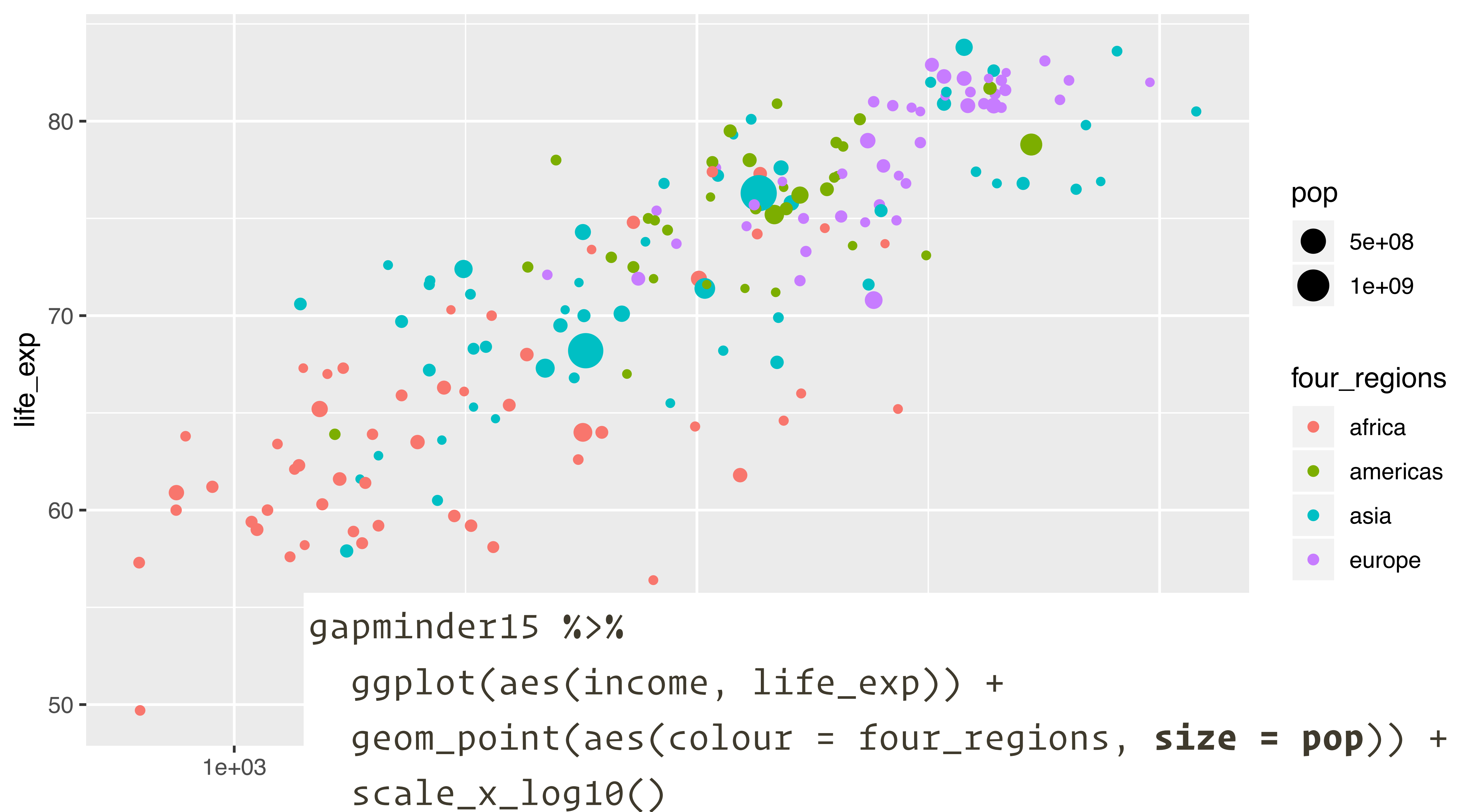












```

# install.packages("devtools")
devtools::install_github("jennybc/frogs")

## Getting to know the frogs

At this point, all we know is that each row is one frog-jump. Frog ids coming ...

```{r}
library(frogs)
library(tidyverse)

frogs
glimpse(frogs)
```

An early figure. Do frogs need to warm up? Do they fatigue? Yes and yes.

```{r frog-fatigue, echo = FALSE}
frogs2 <- frogs %>%
 filter(jump_n < 7) %>%
 mutate(
 jump_n = as.factor(as.integer(jump_n))
)
ggplot(frogs2, aes(x = distance, color = jump_n)) +
 geom_density()
```

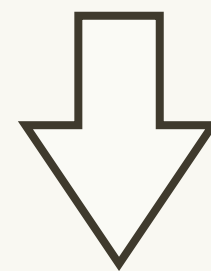
Do professional frog jumping teams get better results? YES.

```{r frog-type, echo = FALSE}
ggplot(frogs, aes(x = distance, color = frog_type)) +
 geom_density()
```

```

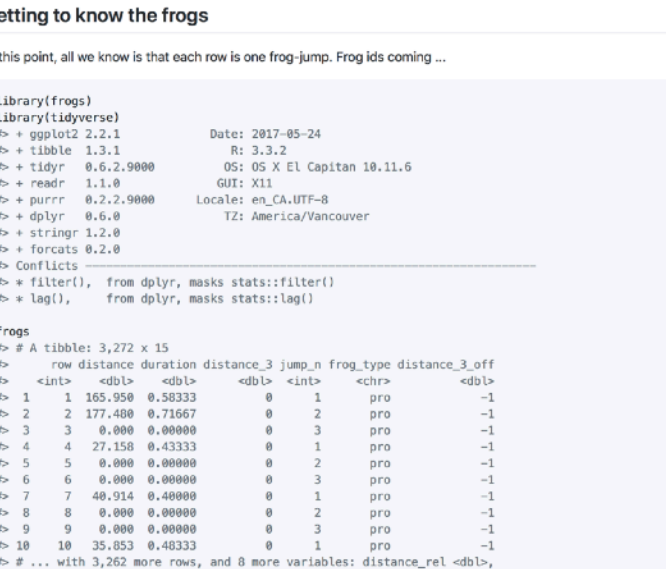
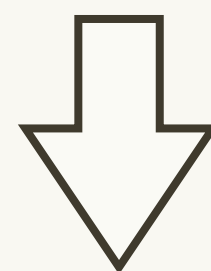
.Rmd

Prose and code



.md

Prose and results



```

library(frogs)
library(tidyverse)
#> ggplot2 2.2.1          Date: 2017-05-24
#> tibble 1.1.1           R: 3.3.2
#> tidyr 0.6.2.9000       OS: OS X El Capitan 10.11.6
#> readr 1.1.0            GUI: X11
#> purrr 0.2.2.9000       Locale: en_CA.UTF-8
#> dplyr 0.6.0             TZ: America/Vancouver
#> string 1.2.0
#> forcats 0.2.0

#> Conflicts —————
#> =========
#> w = filter(), from dplyr, masks stats::filter()
#> = lag(), from dplyr, masks stats::lag()

frogs
#> # A tibble: 3,272 x 15
#>   row distance distance_3_jump_n frog_type distance_3_off
#>   <int> <dbl> <dbl> <dbl> <int> <chr> <dbl>
#> 1 1 165.950 0.50333 0 1 1 pro -1
#> 2 2 177.400 0.71607 0 2 pro -1
#> 3 3 0.000 0.00000 0 3 pro -1
#> 4 4 27.158 0.43333 0 1 pro -1
#> 5 5 0.000 0.00000 0 2 pro -1
#> 6 6 0.000 0.00000 0 3 pro -1
#> 7 7 48.914 0.40000 0 1 pro -1
#> 8 8 0.000 0.00000 0 2 pro -1
#> 9 9 0.000 0.00000 0 3 pro -1
#> 10 10 35.853 0.40333 0 1 pro -1
#> # ... with 3,262 more rows, and 0 more variables: distance_rel_dbl,
#>   day_dbl, angle_at_dbl, angle_10_dbl, angle_80_dbl,
#>   velocity_0_dbl, velocity_10_dbl, velocity_80_dbl
glimpse(frogs)
#> Observations: 3,272
#> Variables: 15

```

.html

.doc

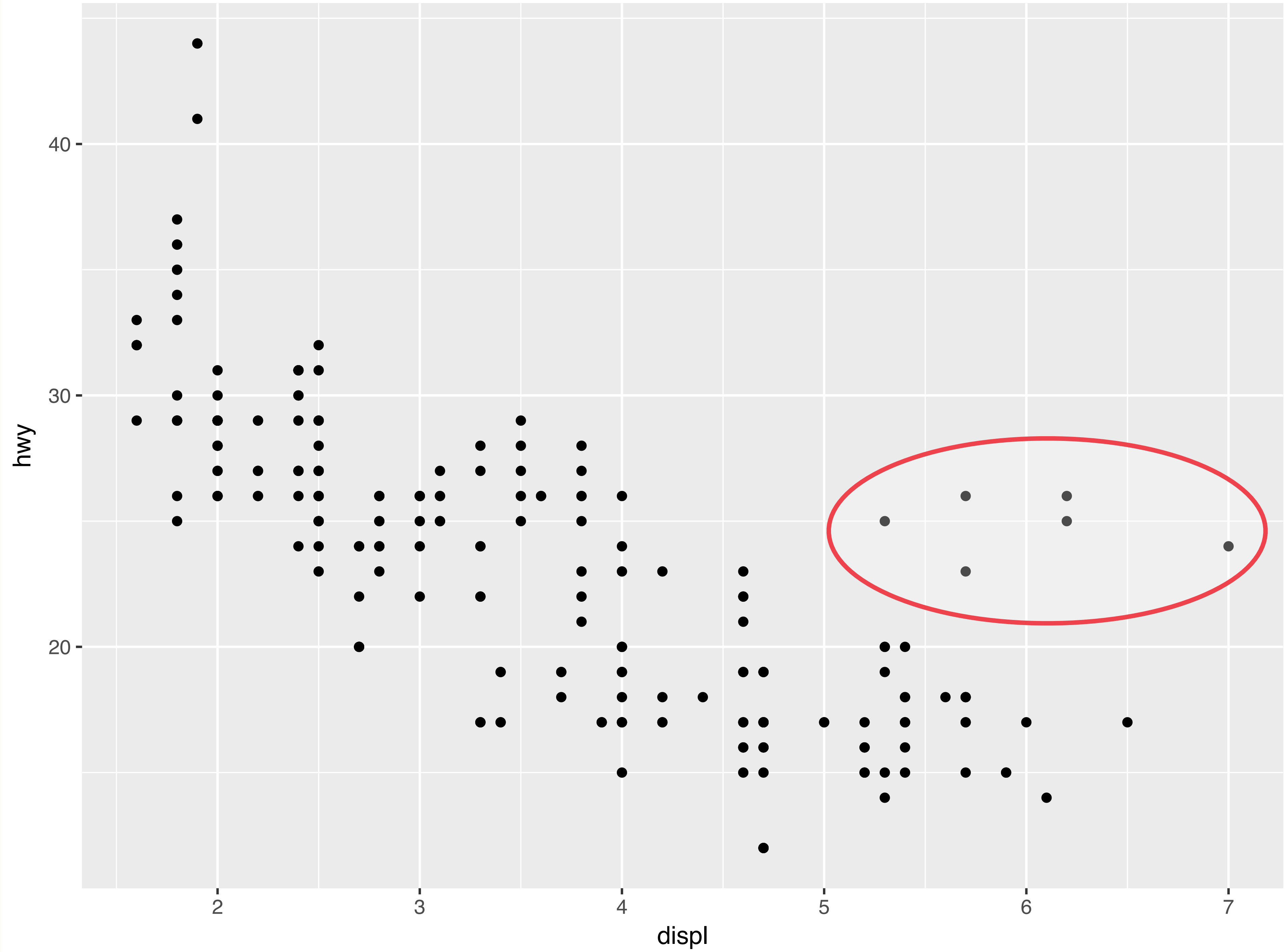
.tex

.ppt

.pdf

● ● ●

But



And this is painful!

```
df %>%  
  rename(  
    date = `Date Created`,  
    name = Name,  
    plays = `Total Plays`,  
    loads = `Total Loads`,  
    apv = `Average Percent Viewed`  
  )
```

Overview

Using Addins

Installation

Running Addins

Keyboard Shortcuts

Developing Addins

Addin Basics

RStudio API

Registering Addins

Execution Modes

Shiny Gadgets

Gadget UI

Gadget Server

Gadget Viewer

Putting It Together

Installation

More Examples

RStudio Addins

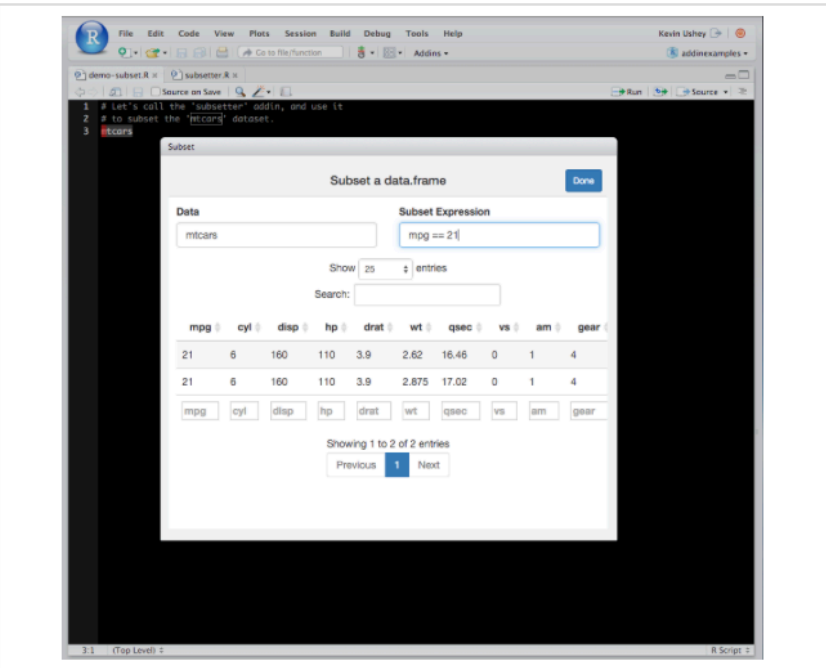
IMPORTANT NOTE: Support for addins is available only within the [most recent release](#) of RStudio (v0.99.878 or later). If you want to try out addins please be sure to download this release.

RStudio addins provide a mechanism for executing R functions interactively from within the RStudio IDE—either through keyboard shortcuts, or through the *Addins* menu.

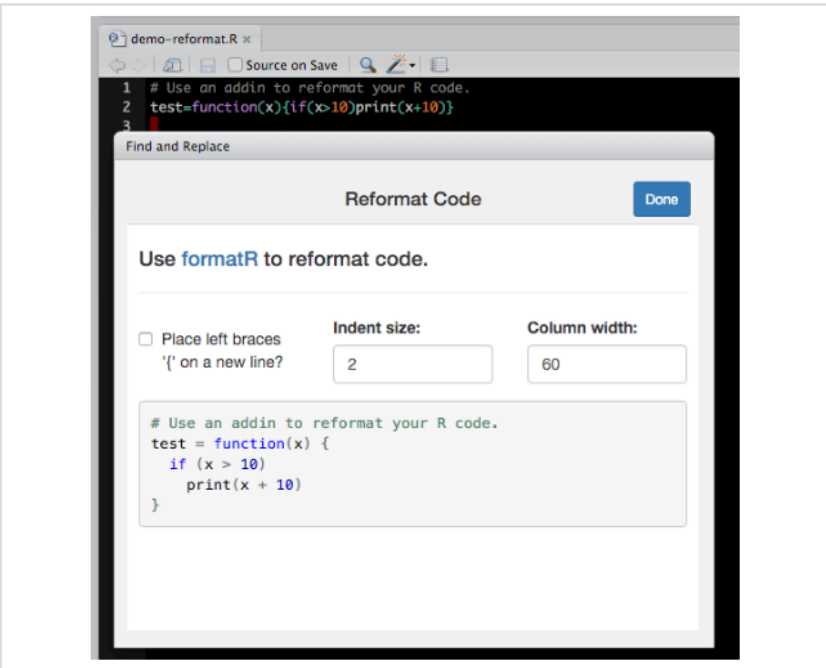
An addin can be as simple as a function that inserts a commonly used snippet of text, and as complex as a Shiny application that accepts input from the user, and later mutates a document open in RStudio. The sky is the limit!

Here are two examples of addins in action (click on the thumbnail to see a brief demonstration):

Subset a dataset



Reformat R code



Using Addins

This guide will walk you through the basics of installing addins, binding keyboard shortcuts to them, and finally developing your own addins.

Installation

RStudio Addins are distributed as [R packages](#). Once you’ve installed an R package that contains addins, they’ll be immediately become available within RStudio.

Let’s start by playing around with a couple of the example addins provided by the [addinexamples](#) package. Within RStudio, install this package (plus its requisite dependencies) with:

```
devtools::install_github("rstudio/addinexamples", type = "source")
```

Running Addins

What next?

```
df %>%
```

```
  filter(n > 1e6) %>%
```

```
  mutate(x = f(y)) %>%
```

```
  ???
```

```
# How predictable is next step from
```

```
# previous steps?
```

Can we do more with autocomplete?

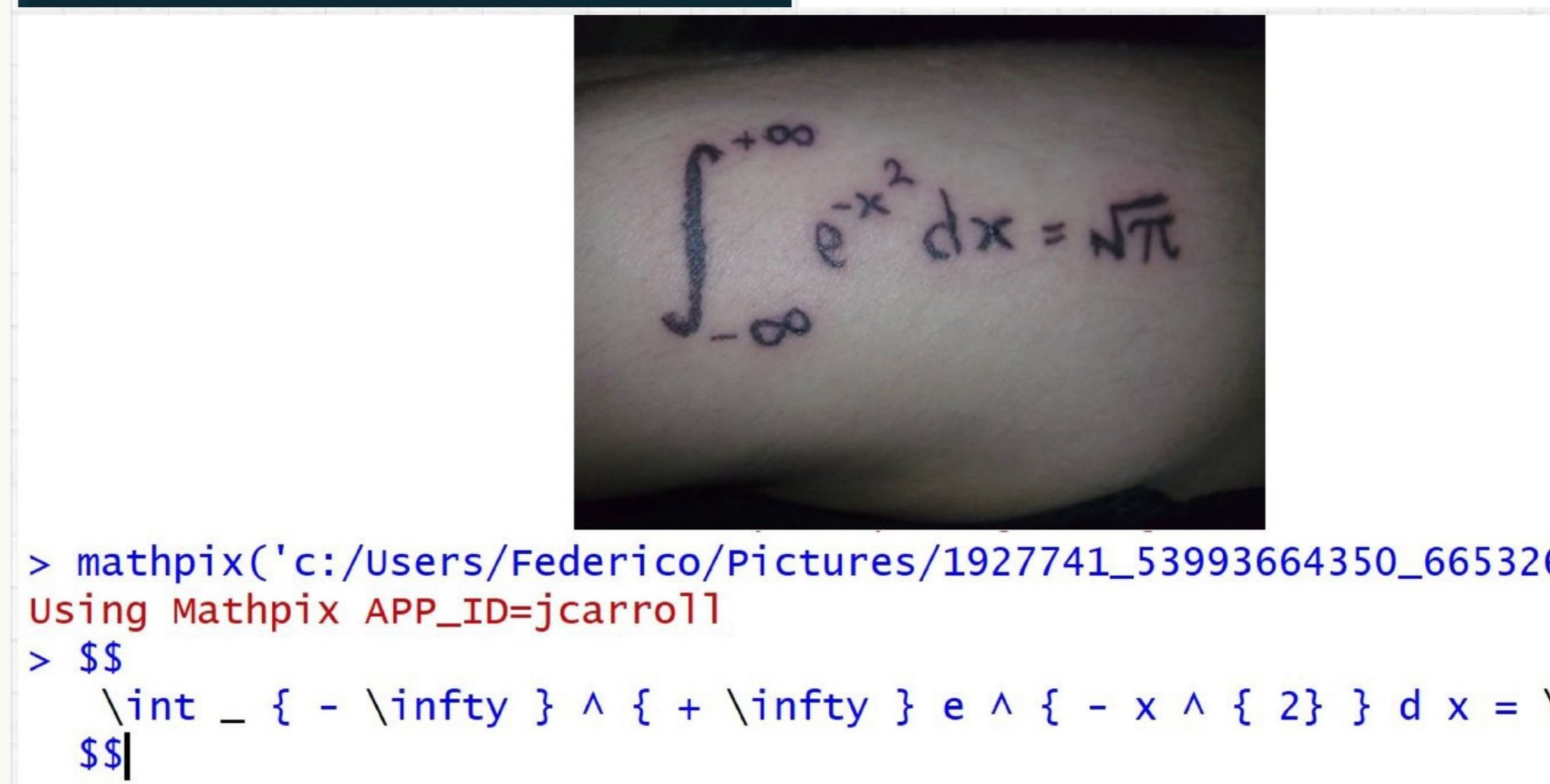
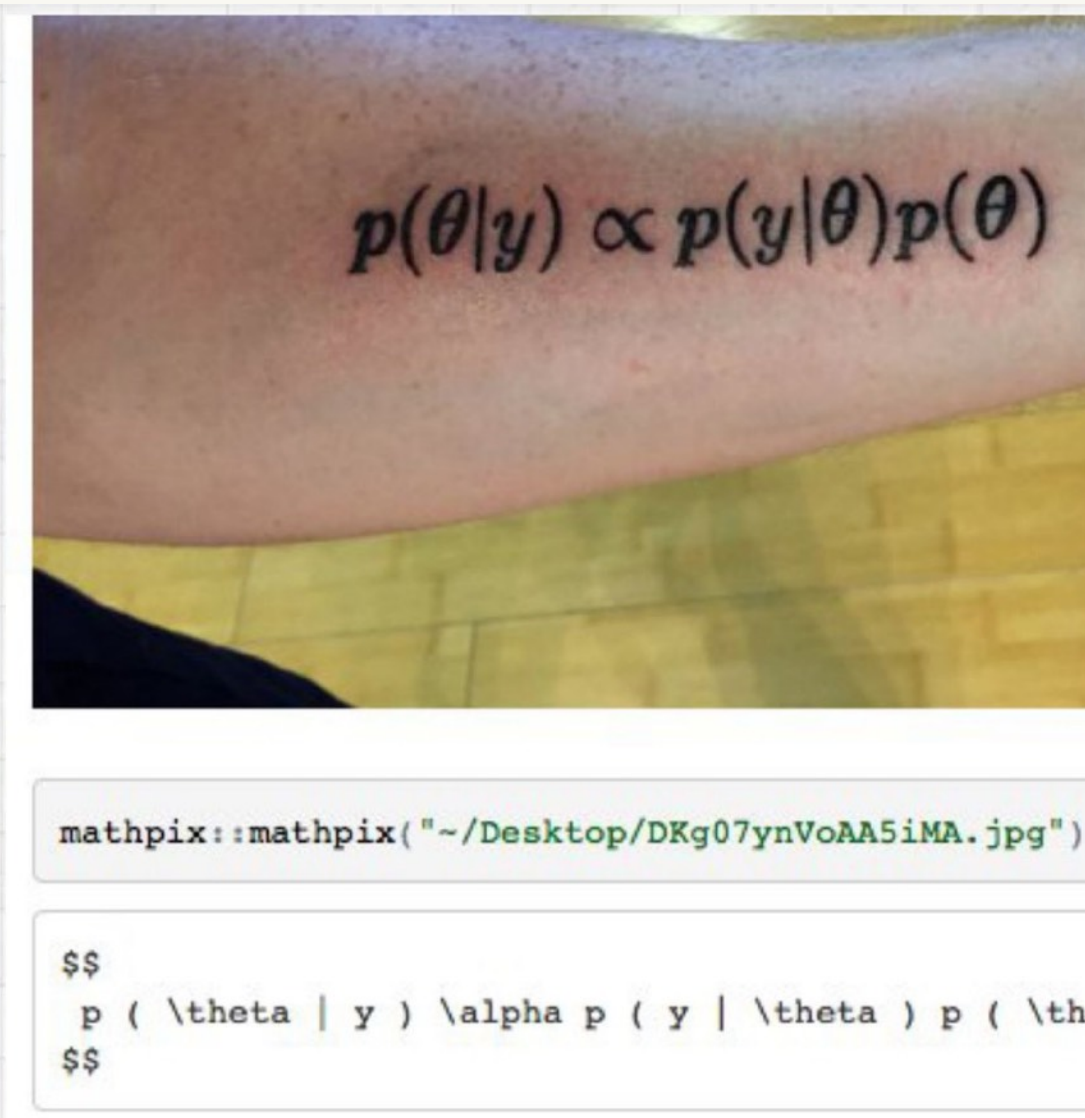
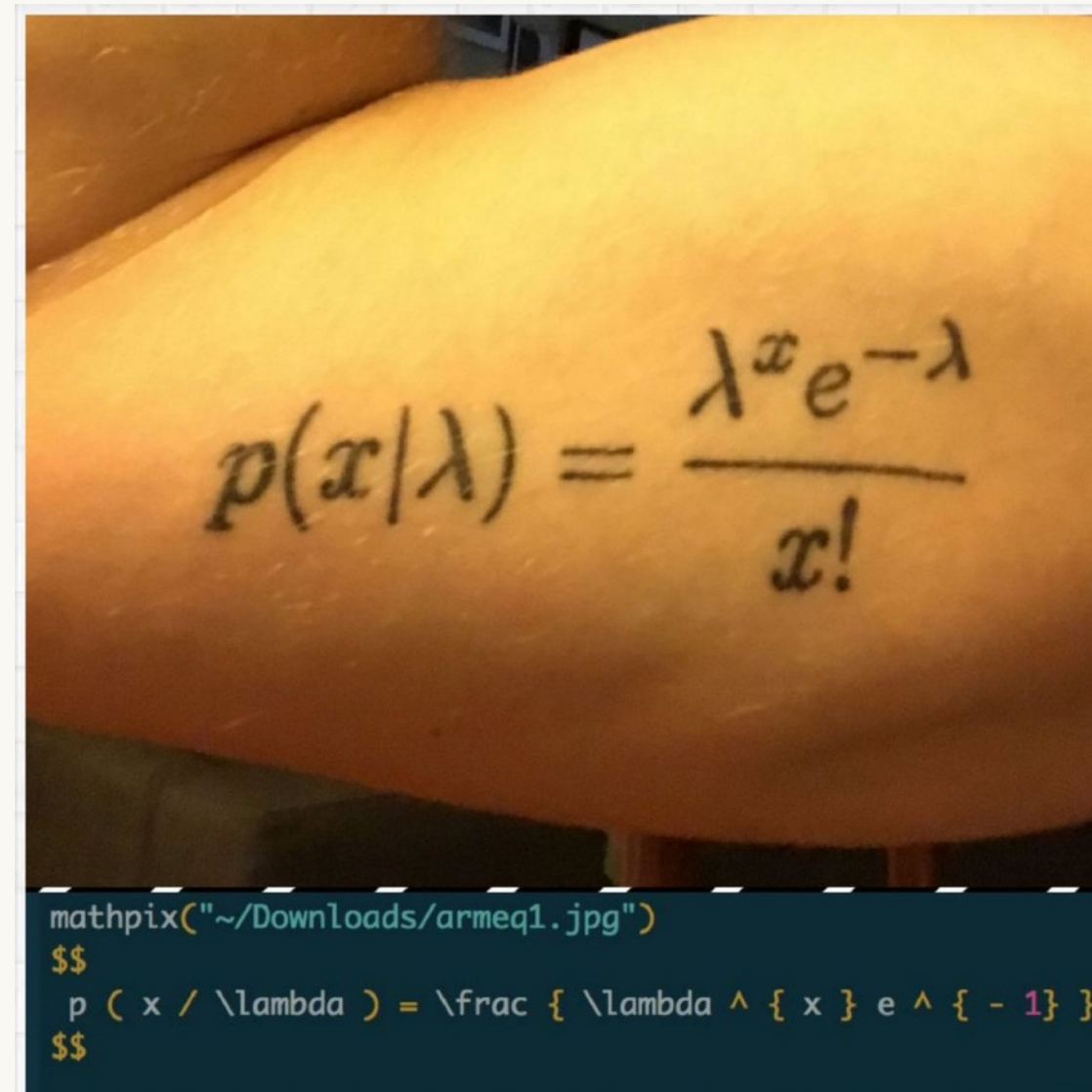
```
>  
>  
>  
>  
>  
>  
>  
>  
>  
> library(  
  p abind  
  p acepack  
  p addcol  
  p ash  
  p assertthat  
  p babynames  
  p backports
```

Where do dialogs and autocomplete intersect?

Learning from examples

- (a) Reported crime in Alabama
- (b) *before:* { 'in', ' ' } 'Alabama' \rightarrow { 'Alabama', *word* }
selection: { 'Alabama' } 'in' \rightarrow { 'in', *word*, *lowercase* }
after: \emptyset ' ' \rightarrow { ' ' }
- (c) *before:* { (' '), ('in', ' '), (*word*, ' '), (*lowercase*, ' ') }
selection: { ('Alabama'), (*word*) }
after: \emptyset
- (d) $\{(), ('Alabama'), ()\}$ $\{(), (~~word~~), ()\}$
 $\{(' '), (), ()\}$ $\{(~~word~~, ' '), (), ()\}$
 $\{(' '), ('Alabama'), ()\}$ $\{(~~word~~, ' '), ('Alabama'), ()\}$
 $\{(~~' '), (~~word~~), ()\}~~$ $\{(~~word~~, ' '), (~~word~~), ()\}$
 $\{('in', ' '), (), ()\}$ $\{(~~lowercase~~, ' '), (), ()\}$
 $\{('in', ' '), ('Alabama'), ()\}$ $\{(~~lowercase~~, ' '), ('Alabama'), ()\}$
 $\{('in', ' '), (~~word~~), ()\}$ $\{(~~lowercase~~, ' '), (~~word~~), ()\}$
- (e) $\{(~~lowercase~~, ' '), ('Alabama'), ()\} \rightarrow /[a-z]+ (Alabama)/$

Figure 10. Regular Expression Inference. (a) The user selects text in a cell. (b) We tokenize selected and surrounding text. For clarity, the figure only includes two neighboring tokens. For each token, we generate a set of matching labels. (c) We enumerate all label sequences matching the text. (d) We then enumerate all candidate *before*, *selection* and *after* combinations. Patterns that do not uniquely match the selected text are filtered (indicated by strike-through). (e) Finally, we construct regular expressions for each candidate pattern.



Fin

We wanted users to be able to begin in an interactive environment, where they did not consciously think of themselves as programming. Then as their needs became clearer and their sophistication increased, they should be able to slide gradually into programming, when the language and system aspects would become more important.

— John Chambers, “Stages in the Evolution of S”



Pit of success

Import

readr
readxl
haven
xml2

Tidy

tibble
tidyr

Transform

dplyr
forcats
hms

lubridate
stringr

Visualise

ggplot2

Model

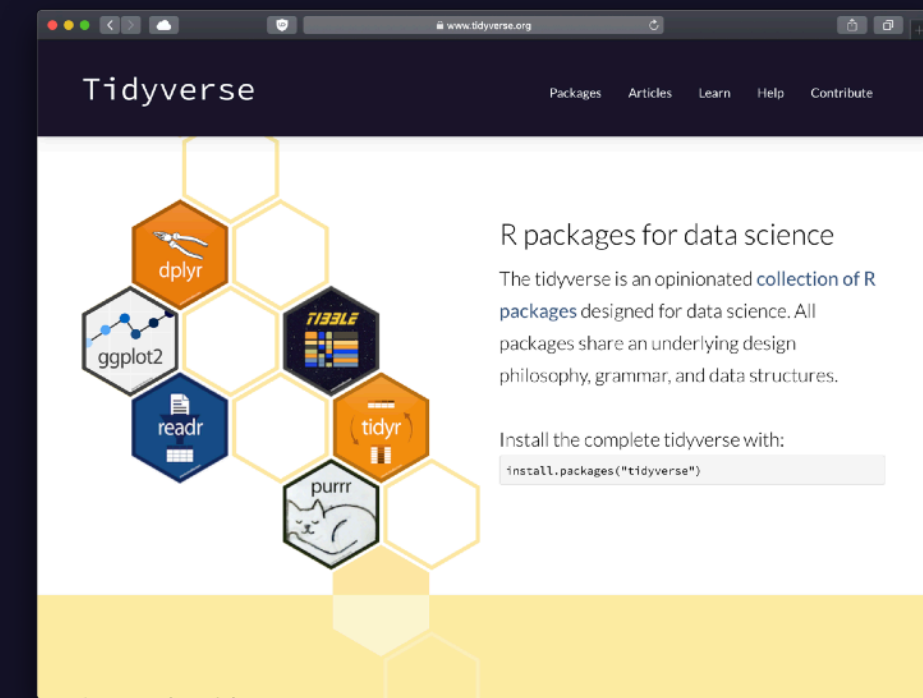
recipes
parsnip

purrr
magrittr

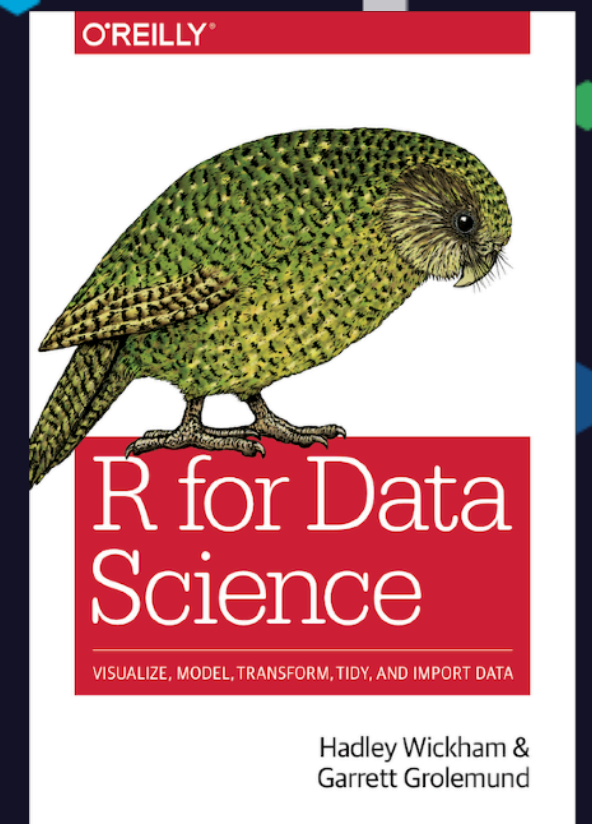
Program

shiny
rmarkdown

Communicate



tidyverse.org



r4ds.had.co.nz

This work is licensed as
Creative Commons
Attribution-ShareAlike 4.0
International

To view a copy of this license, visit
<https://creativecommons.org/licenses/by-sa/4.0/>