# Data Science in R

Dr. Thijs Janzen

```
table1
#> # A tibble: 6 x 4
#>   country     year  cases population
#>   <chr>      <int>  <int>      <int>
#> 1 Afghanistan 1999    745   19987071
#> 2 Afghanistan 2000   2666   20595360
#> 3 Brazil      1999  37737  172006362
#> 4 Brazil      2000  80488  174504898
#> 5 China       1999 212258 1272915272
#> 6 China       2000 213766 1280428583
```
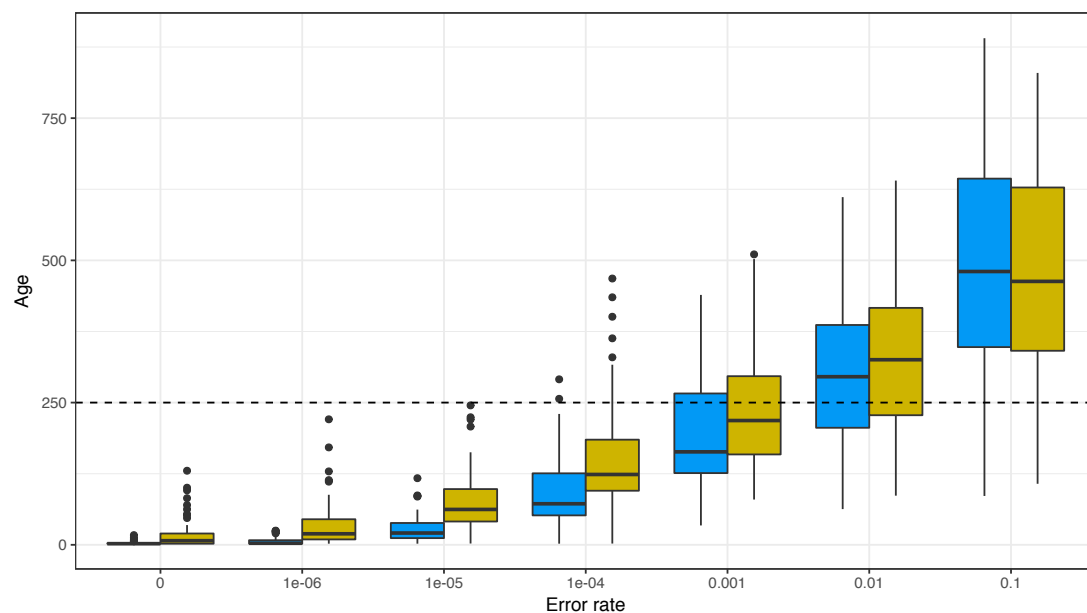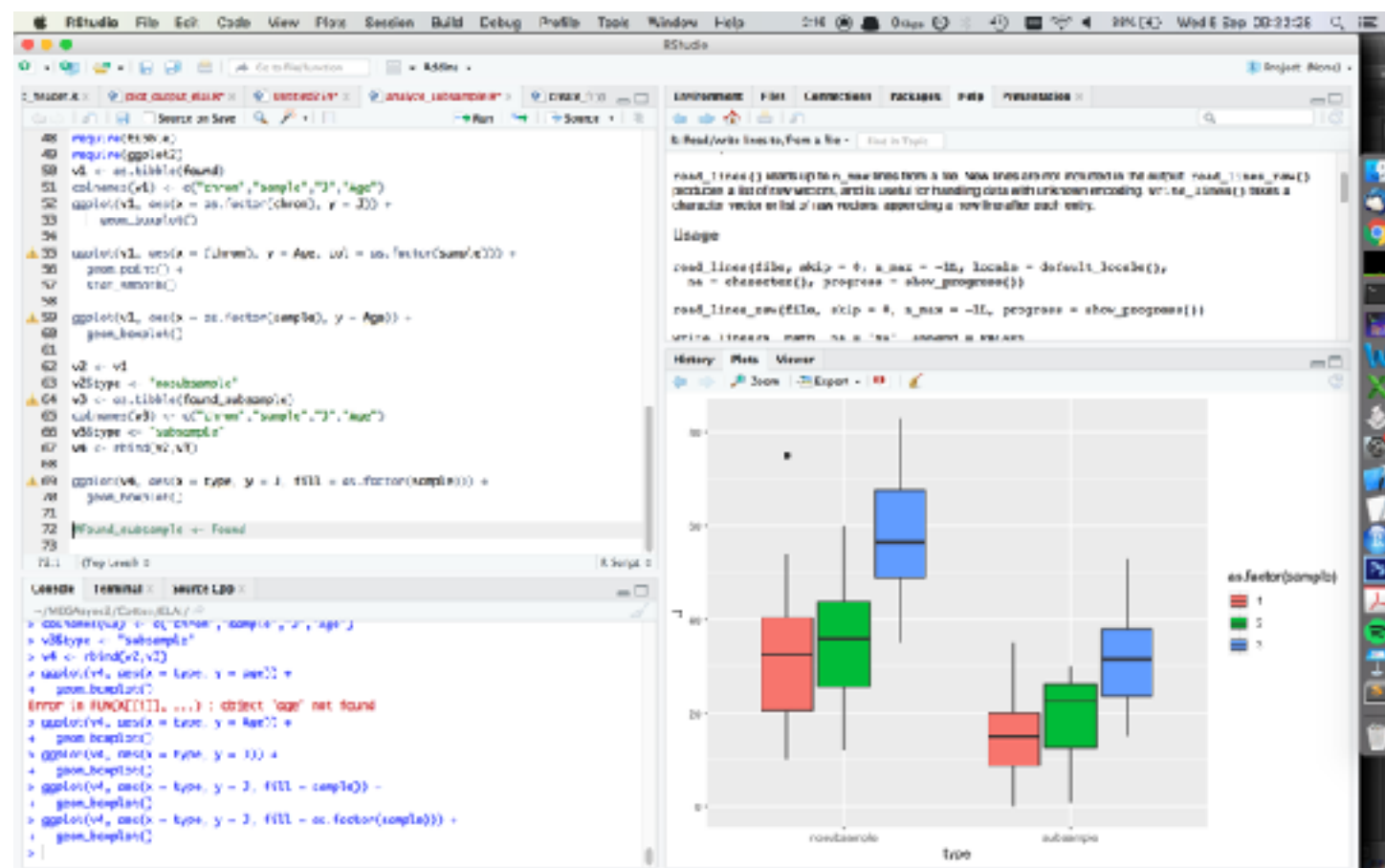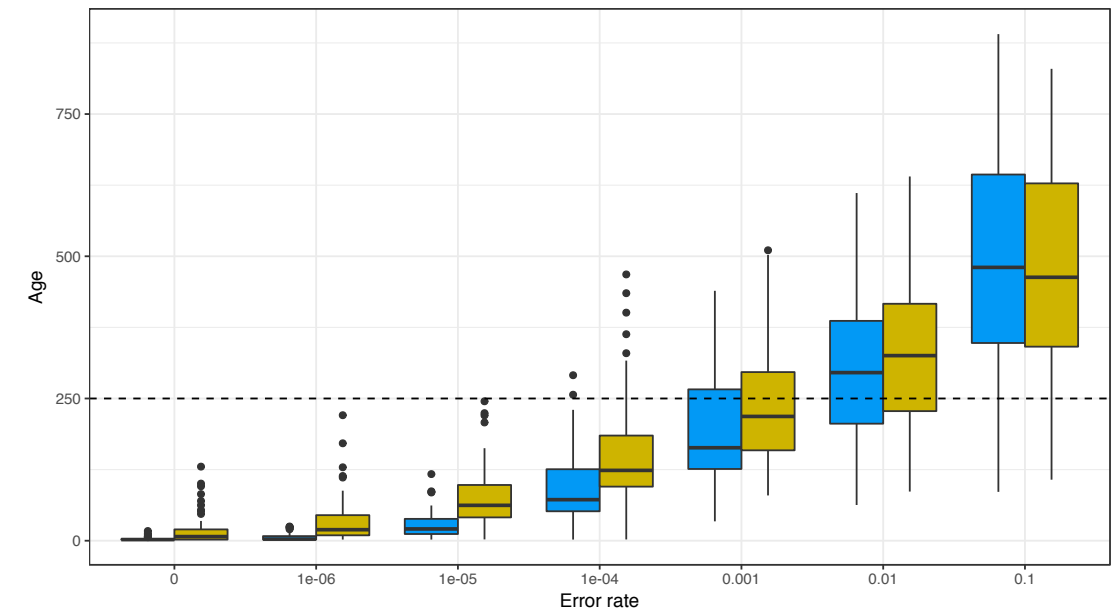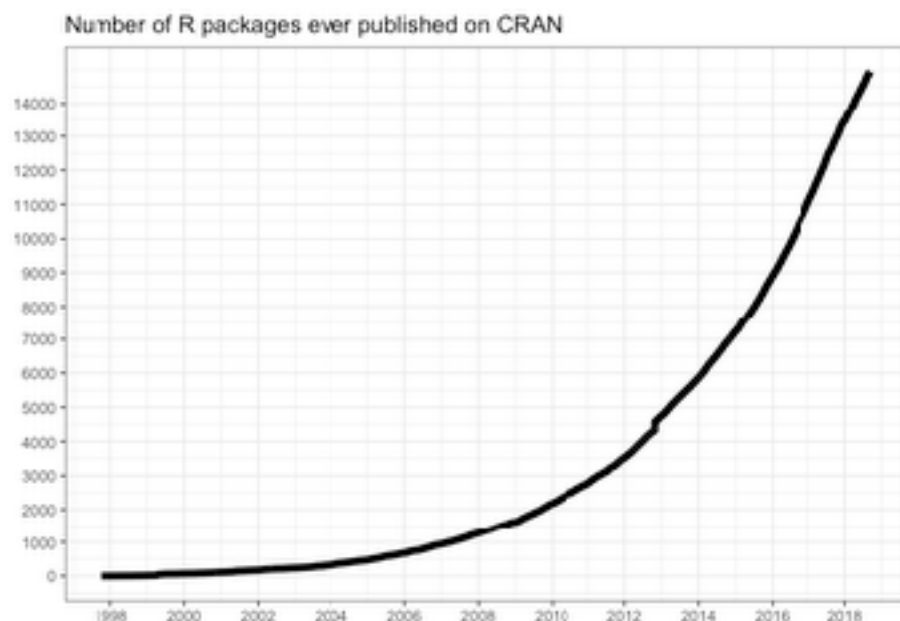
# What is R?

- Environment for statistical computing and graphics

- Free and open source

- Interpreted programming language

# Why use R?

- Statistics

- Data visualisation

- Processing and tidying data

- Reproducible research
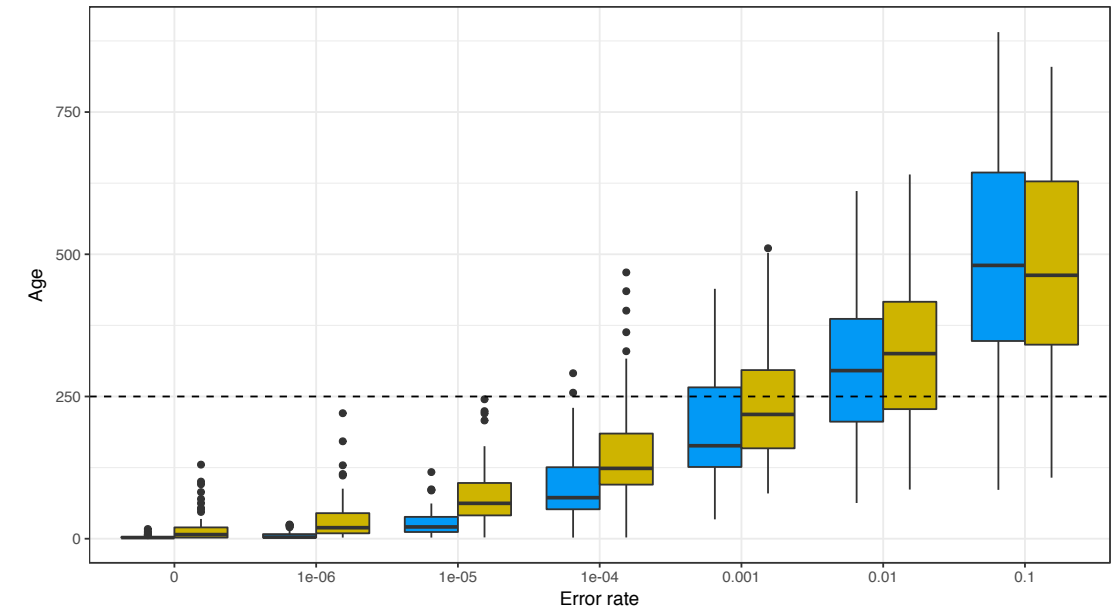
- Many available custom packages

# Why use R?

- Statistics

- **Data visualisation**

- **Processing and tidying data**

- Reproducible research

- Many available custom packages





Peak time of day for sports and leisure

Number of participants throughout the day compared to peak popularity.
Note the morning-and-evening everyday workouts, the midday hobbies,
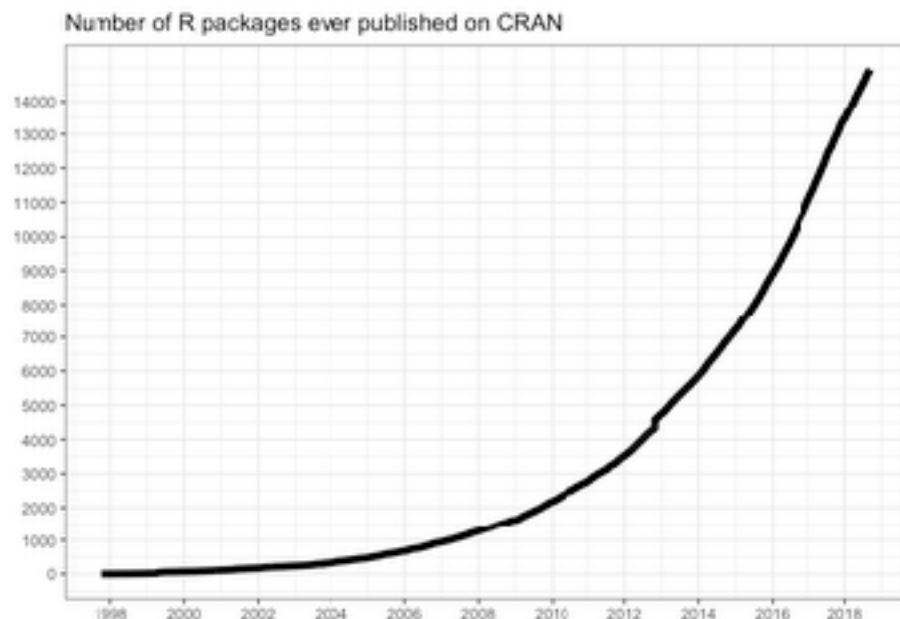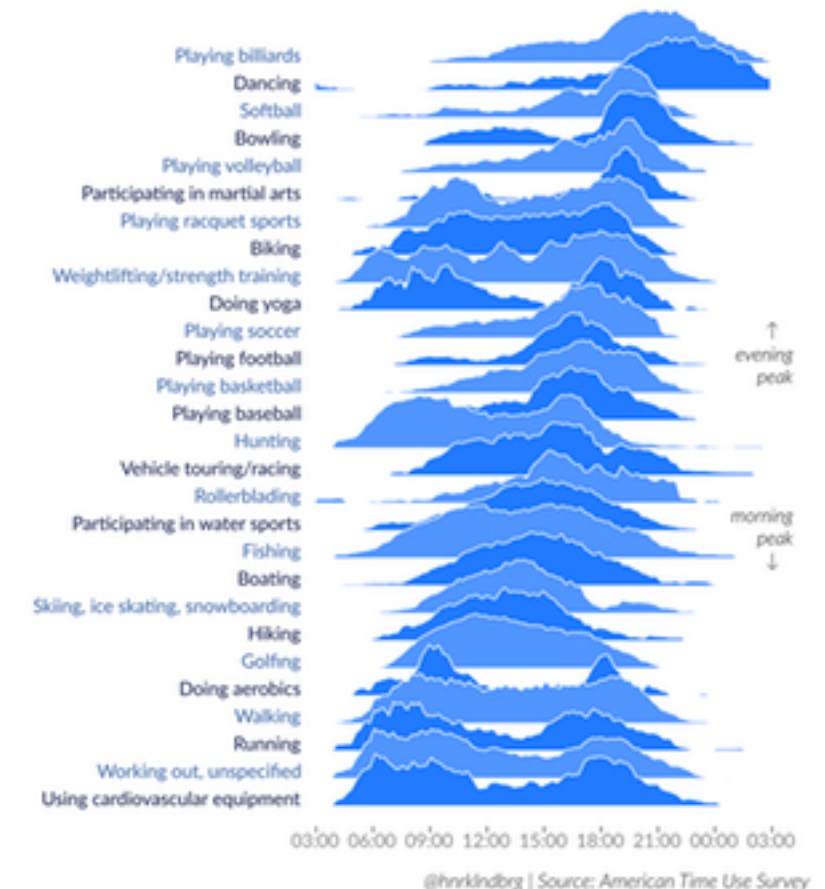and the evenings/late nights out.

@hnrklndbrg | Source: American Time Use Survey



Number of R packages ever published on CRAN

It's easy when you start out programming to get really frustrated and think, "Oh it's me, I'm really stupid," or, "I'm not made out to program." But, that is absolutely not the case. Everyone gets frustrated. I still get frustrated occasionally when writing R code. It's just a natural part of programming. So, it happens to everyone and gets less and less over time. Don't blame yourself. Just take a break, do something fun, and then come back and try again later.

**Hadley Wickham,**
**Chief Scientist at Rstudio**
**Developer of the tidyverse**

# Goal of today

- Load data into R

- Re-structure data to improve handling: 'tidying data'

- Plot results

- Understand basics of tidy workflow



Species richness across sand cover

# Requirements

- Rstudio (no strict requirement, but makes life easy)

- tidyverse packages:

  - tibble

  - readr

  - ggplot2

  - dplyr

```r
install.packages("tidyverse")
library(tidyverse)
```

# Where to find files for this workshop

https://github.com/thijsjanzen/youmares_workshop_R

# tidyverse

# Data structure

# Tidy data

- Each variable has it's own column

- Each observation has it's own row

- Each value has it's own cell

# Examples tidy data

```
table1

#> # A tibble: 6 x 4
#>   country      year  cases population
#>   <chr>       <int>  <int>      <int>
#> 1 Afghanistan  1999    745   19987071
#> 2 Afghanistan  2000   2666   20595360
#> 3 Brazil       1999  37737  172006362
#> 4 Brazil       2000  80488  174504898
#> 5 China        1999 212258 1272915272
#> 6 China        2000 213766 1280428583
```

| Movie | Race | Sex | Words |
|-------|------|-----|-------|
| <chr> | <chr> | <chr> | <chr> |
| 1 Fellowship of the Ring | Elf | Female | 1229 |
| 2 Fellowship of the Ring | Hobbit | Female | 14 |
| 3 Fellowship of the Ring | Man | Female | 0 |
| 4 Fellowship of the Ring | Elf | Male | 971 |
| 5 Fellowship of the Ring | Hobbit | Male | 3644 |
| 6 Fellowship of the Ring | Man | Male | 1995 |
| 7 Two towers | Elf | Female | 331 |
| 8 Two towers | Hobbit | Female | 0 |
| 9 Two towers | Man | Female | 401 |
| 10 Two towers | Elf | Male | 513 |
| 11 Two towers | Hobbit | Male | 2463 |
| 12 Two towers | Man | Male | 3589 |

# But... my data is not tidy?

- When recording data, your data is often not tidy

- There are two functions (amongst others) to help you make your data tidy:

  - gather

  - spread

# Importing data into R

read_tsv, read_csv, read_delim

read_tsv(file, col_names = TRUE)

read_csv(file, col_names = TRUE)

read_delim(file, delim, col_names = TRUE)

# Reading data into R

```
fish_counts <- read_tsv(file = "cichlid_plots.txt")

fish_counts


lotr_words <- read_tsv(file = "lotr_words.txt")

lotr_words
```

# Tidying data: gather

- The function `gather` combines multiple columns into one column, and adds an extra indicator column.

- gather(data, key, value, columns)

  - `data` = the data to be converted

  - `key` = variable name that is going to contain the column name

  - `value` = variable name that is going to contain the gathered data

  - `columns` = selection of which columns need to be gathered

# Example gather

```
fellow <- read_tsv("fellowship.txt")
```

```
> fellow
# A tibble: 3 x 3
  Race    Female   Male
  <chr>    <int>  <int>
1 Elf       1229    971
2 Hobbit      14   3644
3 Man          0   1995
```

```
fellow_gathered <- gather(data    = fellow,
                          key      =
                          value    =
                          columns  =
```

# Example gather

```
fellow <- read_tsv("fellowship.txt")
```

```
> fellow
# A tibble: 3 x 3
  Race    Female  Male     key
  <chr>    <int> <int>
1 Elf       1229   971
2 Hobbit      14  3644
3 Man          0  1995
```

```
fellow_gathered <- gather(data     = fellow,
                          key      =
                          value    =
                          columns  =
```

# Example gather

```
fellow <- read_tsv("fellowship.txt")
```

```
> fellow
# A tibble: 3 x 3
  Race    Female  Male     key
  <chr>    <int> <int>
1 Elf       1229   971
2 Hobbit      14  3644
3 Man          0  1995
```

```
fellow_gathered <- gather(data    = fellow,
                          key     = "Sex",
                          value   =
                          columns =
```

# Example gather

```
fellow <- read_tsv("fellowship.txt")
```

```
> fellow
# A tibble: 3 x 3
    Race    Female   Male      key
    <chr>   <int>  <int>
1 Elf       1229     971
2 Hobbit      14    3644       value
3 Man          0    1995
```

```
fellow_gathered <- gather(data    = fellow,
                          key      = "Sex",
                          value    =
                          columns  =
```

# Example gather

```
fellow <- read_tsv("fellowship.txt")
```

```
> fellow
# A tibble: 3 x 3
  Race    Female  Male        key
  <chr>   <int>   <int>
1 Elf     1229    971
2 Hobbit  14      3644        value
3 Man     0       1995
```

```
fellow_gathered <- gather(data    = fellow,
                          key     = "Sex",
                          value   = "Words",
                          columns =
```

# Example gather

```
fellow <- read_tsv("fellowship.txt")
```

```
> fellow
# A tibble: 3 x 3
  Race    Female  Male     key
  <chr>    <int> <int>
1 Elf       1229   971
2 Hobbit      14  3644     value
3 Man          0  1995
```

```
fellow_gathered <- gather(data    = fellow,
                          key     = "Sex",
                          value   = "Words",
                          columns = c("Male","Female"))
```

# Example gather

```
fellow <- read_tsv("fellowship.txt")
```

```
> fellow
# A tibble: 3 x 3
  Race    Female   Male      key
  <chr>   <int>   <int>
1 Elf      1229     971
2 Hobbit     14    3644      value
3 Man         0    1995
```

```
fellow_gathered <- gather(data    = fellow,
                          key     = "Sex",
                          value   = "Words",
                          columns = c("Male","Female"))
```

```
> fellow_gathered
# A tibble: 6 x 3
  Race    Sex      Words
  <chr>   <chr>    <int>
1 Elf     Male       971
2 Hobbit  Male      3644
3 Man     Male      1995
4 Elf     Female    1229
5 Hobbit  Female      14
6 Man     Female       0
```

# How to indicate the columns?

- Use the names:

```
gather(fellow, key = "Sex", value = "Words", "Female", "Male")
```

- Use the index:

```
gather(fellow, key = "Sex", value = "Words", 2:3)
```

- Use all columns (except the first):

```
gather(fellow, key = "Sex", value = "Words", -1)
```

# Plotting data

- OK, we have tidy data now

- How to visualise results?

# ggplot

- ggplot: the Grammar of Graphics

- Plots are constructed out of building blocks:

  - data

  - aesthetic mapping

  - geometric object

  - statistical transformations

  - scales

  - coordinate systems

  - labels

# ggplot

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

    geom_        +

    stat_        +

    xlab()       +
```

**aesthetics**: indicate what is on the x axis, on the y-axis,
            and if you need grouping of your data

**geom_point / geom_line / geom_bar** etc. : indicates the type of plot
                        (scatter, line, barplot, box plot etc)

**stat_smooth()** : indicates additional statistics

# plotting lotr

**Let's create a bar plot, split per race and sex**

```
ggplot(data = fellow_gathered, aes(x = Race, y = Words, fill = Sex)) +
        geom_bar(stat = "identity", position = "dodge")
```

# plotting lotr

```
> fellow_gathered
# A tibble: 6 x 3
   Race   Sex     Words
   <chr>  <chr>   <int>
1  Elf    Male      971
2  Hobbit Male     3644
3  Man    Male     1995
4  Elf    Female   1229
5  Hobbit Female     14
6  Man    Female      0
```

**Let's create a bar plot, split per race and sex**

```
ggplot(data = fellow_gathered, aes(x = Race, y = Words, fill = Sex)) +
        geom_bar(stat = "identity", position = "dodge")
```

# plotting lotr

```
> fellow_gathered
# A tibble: 6 x 3
  Race   Sex     Words
  <chr>  <chr>   <int>
1 Elf    Male      971
2 Hobbit Male     3644
3 Man    Male     1995
4 Elf    Female   1229
5 Hobbit Female     14
6 Man    Female      0
```

**Let's create a bar plot, split per race and sex**

```
ggplot(data = fellow_gathered, aes(x = Race, y = Words, fill = Sex)) +
        geom_bar(stat = "identity", position = "dodge")
```

# ggplot

- `ggplot(fish_counts, aes(x = sand_cover, y = number_of_species)) +`
  `geom_point()`

  - aesthetic mapping

  - geometric object

# Examples gather

**The fish counts dataset contains one column for every species, but rather, we would like to have one column indicating species, and one column indicating the observed number.**

**How can we do that with gather?**

```
fish_counts2 <- gather(data     = fish_counts,
                       key      =
                       value    =
                       columns  =
```

# Examples gather

**The fish counts dataset contains one column for every species, but rather, we would like to have one column indicating species, and one column indicating the observed number.**

**How can we do that with gather?**

```
fish_counts2 <- gather(data    = fish_counts,
                       key     = "Species",
                       value   =
                       columns =
```

# Examples gather

**The fish counts dataset contains one column for every species, but rather, we would like to have one column indicating species, and one column indicating the observed number.**

**How can we do that with gather?**

```
fish_counts2 <- gather(data    = fish_counts,
                       key     = "Species",
                       value   = "Count",
                       columns =
```

# Examples gather

**The fish counts dataset contains one column for every species, but rather, we would like to have one column indicating species, and one column indicating the observed number.**

**How can we do that with gather?**

```
fish_counts2 <- gather(data    = fish_counts,
                       key     = "Species",
                       value   = "Count",
                       columns = -c(1:6) )
```

# Examples gather

The fish counts dataset contains one column for every species, but rather, we would like to have one column indicating species, and one column indicating the observed number.

How can we do that with gather?

```
fish_counts2 <- gather(data    = fish_counts,
                       key     = "Species",
                       value   = "Count",
                       columns = -c(1:6) )
```

```
> fish_counts2
# A tibble: 1,764 x 8
    Plot number_of_indivi… number_of_speci… sand_cover depth rugosity Species Count
   <int>            <int>            <int>      <dbl> <dbl>    <dbl> <chr>   <int>
1     1              135               17       13.4  13.8     1.55 Altolam…    1
2     2              217               16       54.4  13.8     1.33 Altolam…    0
3     3              172               24        6.98 13.0     1.46 Altolam…    2
4     4               74               21       19.7  15.1     1.15 Altolam…    0
5     5               79               21       38.1  14.4     1.47 Altolam…    0
6     6               65               16       75.4  12.8     1.61 Altolam…    0
7     7              338               26        7.77 11.0     1.34 Altolam…    6
8     8              446               25        4.84  9.8     1.27 Altolam…    0
9     9              310               26       14.4   8.15    1.31 Altolam…    2
10    10             577               28        0    11.2     1.70 Altolam…    2
# … with 1,754 more rows
```

# Spreading

- Collecting observations that are scattered among multiple rows

- Spreading is the exact opposite of gather

```
> table2
# A tibble: 12 x 4
           country  year         type      count
             <chr> <int>        <chr>      <int>
 1 Afghanistan  1999        cases        745
 2 Afghanistan  1999 population   19987071
 3 Afghanistan  2000        cases       2666
 4 Afghanistan  2000 population   20595360
 5       Brazil  1999        cases      37737
 6       Brazil  1999 population  172006362
 7       Brazil  2000        cases      80488
 8       Brazil  2000 population  174504898
 9        China  1999        cases     212258
10        China  1999 population 1272915272
11        China  2000        cases     213766
12        China  2000 population 1280428583
```

- spread(data, key, value)

  - data = the data to be converted

  - key = variable name that needs to be spread

  - value = column of corresponding data values

```
table2_b <- spread(table2, key = "type", value = "count")
```

```
# A tibble: 6 x 4
        country  year   cases population
*         <chr> <int>   <int>      <int>
1 Afghanistan   1999     745   19987071
2 Afghanistan   2000    2666   20595360
3      Brazil   1999   37737  172006362
4      Brazil   2000   80488  174504898
5       China   1999  212258 1272915272
6       China   2000  213766 1280428583
```

# Gather

# Spread



| country | year | key | value |
|---|---|---|---|
| Afghanistan | 1999 | cases | 745 |
| Afghanistan | 1999 | population | 19987071 |
| Afghanistan | 2000 | cases | 2666 |
| Afghanistan | 2000 | population | 20595360 |
| Brazil | 1999 | cases | 37737 |
| Brazil | 1999 | population | 172006362 |
| Brazil | 2000 | cases | 80488 |
| Brazil | 2000 | population | 174504898 |
| China | 1999 | cases | 212258 |
| China | 1999 | population | 1272915272 |
| China | 2000 | cases | 213766 |
| China | 2000 | population | 1280428583 |

table2

| country | year | cases | population |
|---|---|---|---|
| Afghanistan | 1999 | 745 | 19987071 |
| Afghanistan | 2000 | 2666 | 20595360 |
| Brazil | 1999 | 37737 | 172006362 |
| Brazil | 2000 | 80488 | 174504898 |
| China | 1999 | 212258 | 1272915272 |
| China | 2000 | 213766 | 1280428583 |

# Gather

| country | year | cases |
|---|---|---|
| Afghanistan | 1999 | 745 |
| Afghanistan | 2000 | 2666 |
| Brazil | 1999 | 37737 |
| Brazil | 2000 | 80488 |
| China | 1999 | 212258 |
| China | 2000 | 213766 |

| country | 1999 | 2000 |
|---|---|---|
| Afghanistan | 745 | 2666 |
| Brazil | 37737 | 80488 |
| China | 212258 | 213766 |

table4

# Spread

| country | year | key | value |
|---|---|---|---|
| Afghanistan | 1999 | cases | 745 |
| Afghanistan | 1999 | population | 19987071 |
| Afghanistan | 2000 | cases | 2666 |
| Afghanistan | 2000 | population | 20595360 |
| Brazil | 1999 | cases | 37737 |
| Brazil | 1999 | population | 172006362 |
| Brazil | 2000 | cases | 80488 |
| Brazil | 2000 | population | 174504898 |
| China | 1999 | cases | 212258 |
| China | 1999 | population | 1272915272 |
| China | 2000 | cases | 213766 |
| China | 2000 | population | 1280428583 |

table2

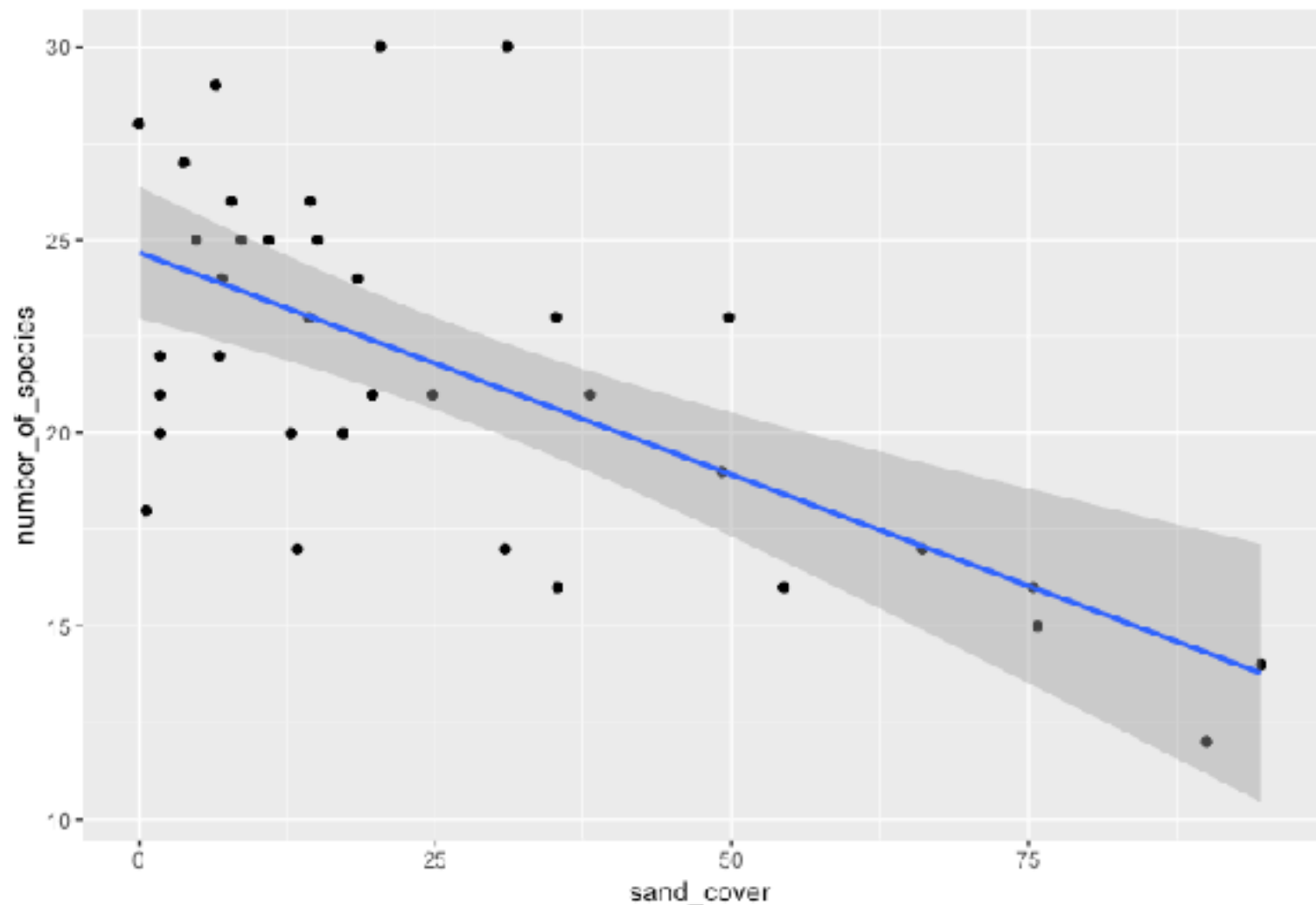| country | year | cases | population |
|---|---|---|---|
| Afghanistan | 1999 | 745 | 19987071 |
| Afghanistan | 2000 | 2666 | 20595360 |
| Brazil | 1999 | 37737 | 172006362 |
| Brazil | 2000 | 80488 | 174504898 |
| China | 1999 | 212258 | 1272915272 |
| China | 2000 | 213766 | 1280428583 |

# ggplot

- ggplot(fish_counts, aes(x = sand_cover, y = number_of_species)) +
            geom_point()

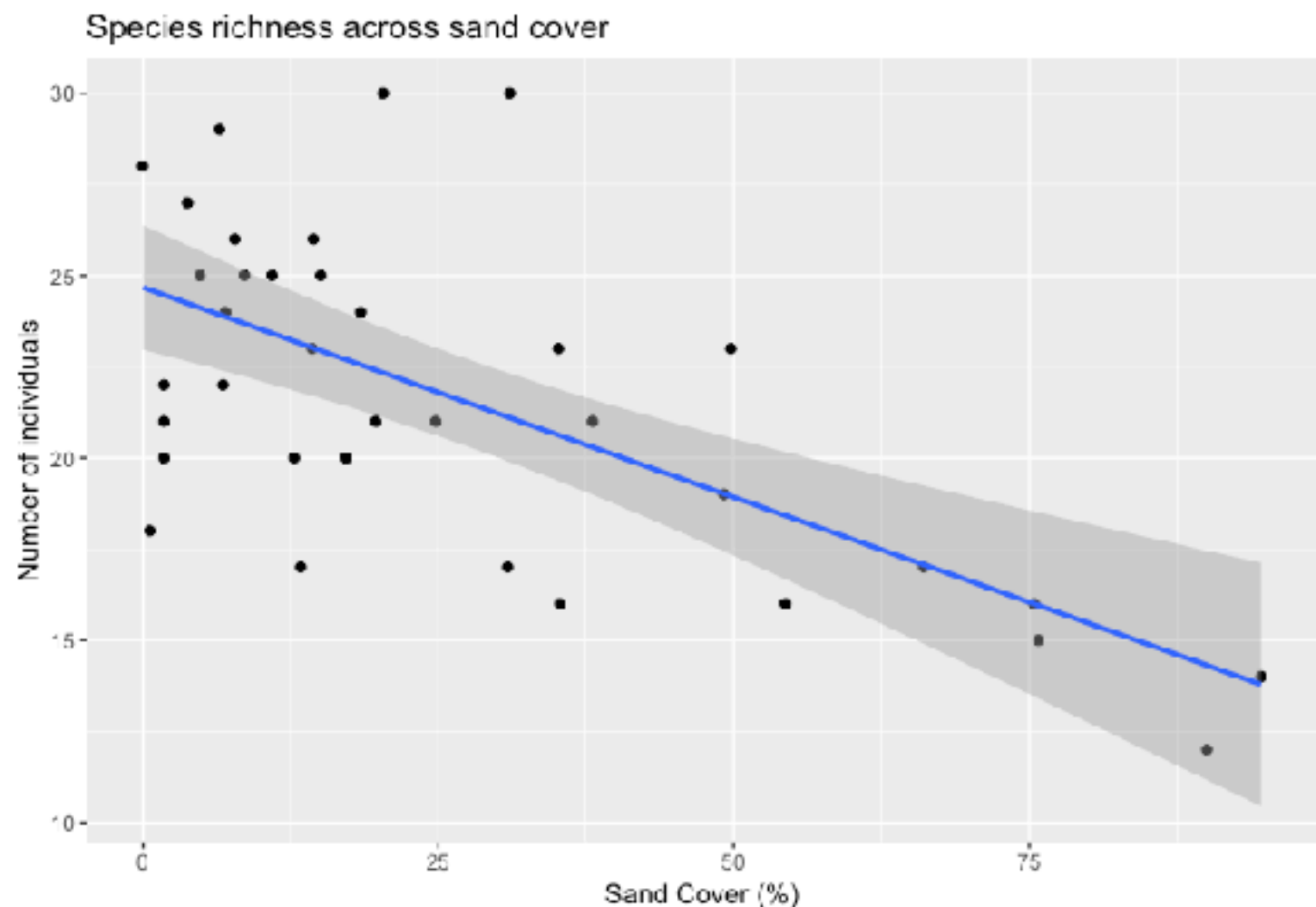  - aesthetic mapping

  - geometric object

# ggplot

- `ggplot(fish_counts, aes(x = sand_cover, y = number_of_species)) +`
  `geom_point() +`
  `stat_smooth(method = "lm")`

- aesthetic mapping

- geometric object

- statistical transformation

# ggplot

- `ggplot(fish_counts, aes(x = sand_cover, y = number_of_species)) +`
    `geom_point() +`
    `stat_smooth() +`
    `xlab("Sand Cover (%)") +`
    `ylab("Number of individuals") +`
    `ggtitle("Species richness across sand cover")`

- aesthetic mapping

- geometric object

- statistical transformation

- labels

# Summary

- Load data into R using `read_tsv`

- Tidy your data using `gather` & `spread`

- Visualize your results using `ggplot`

# Exercises

- If you have your own data:

    - is it tidy?

    - what do you need to do to make it tidy?


-

# Thank you!

- Further reading:

  - R for data science
    (free on http://r4ds.had.co.nz/)

  - ggplot:
    https://tutorials.iq.harvard.edu/R/
    Rgraphics/Rgraphics.html