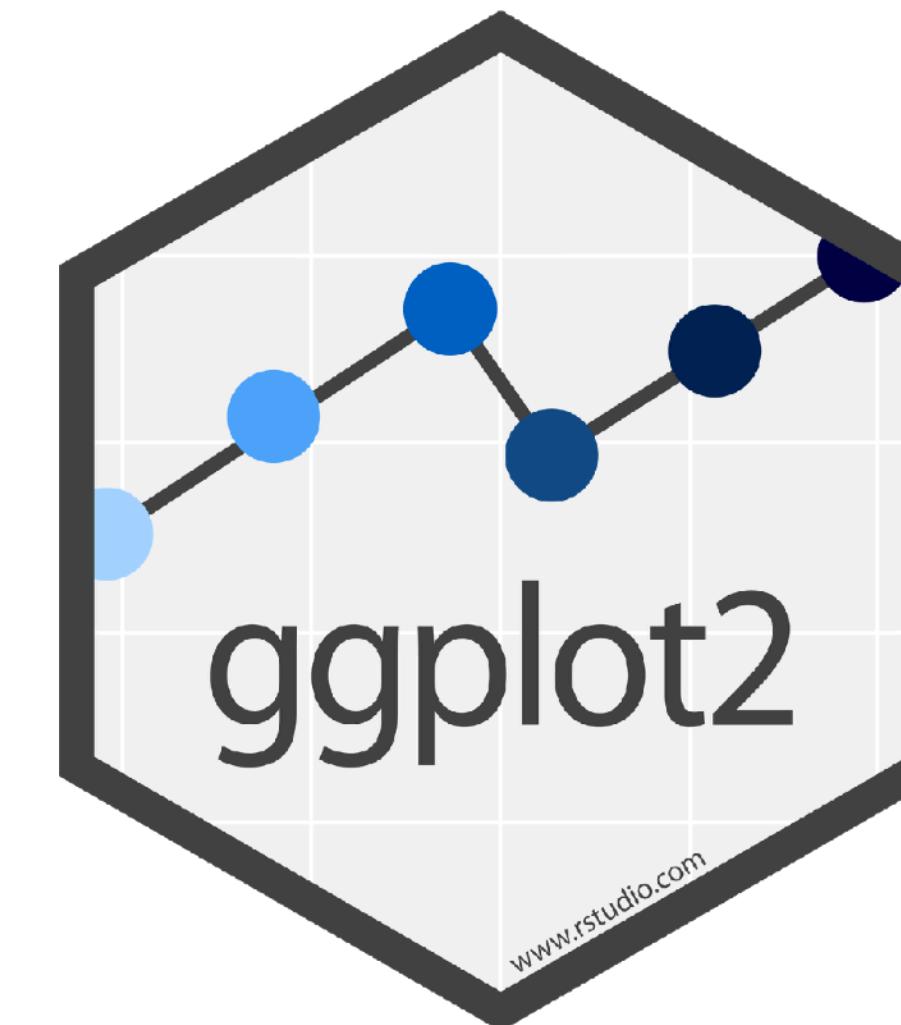


# Programme for today

- Introductory lecture
- Hands-on dataviz tutorial
- Individual exercises



At the end of the day you should be familiar with making visualisations in R using the ggplot2 package



The greatest value of a picture  
is when it forces us to notice  
what we never expected to see.

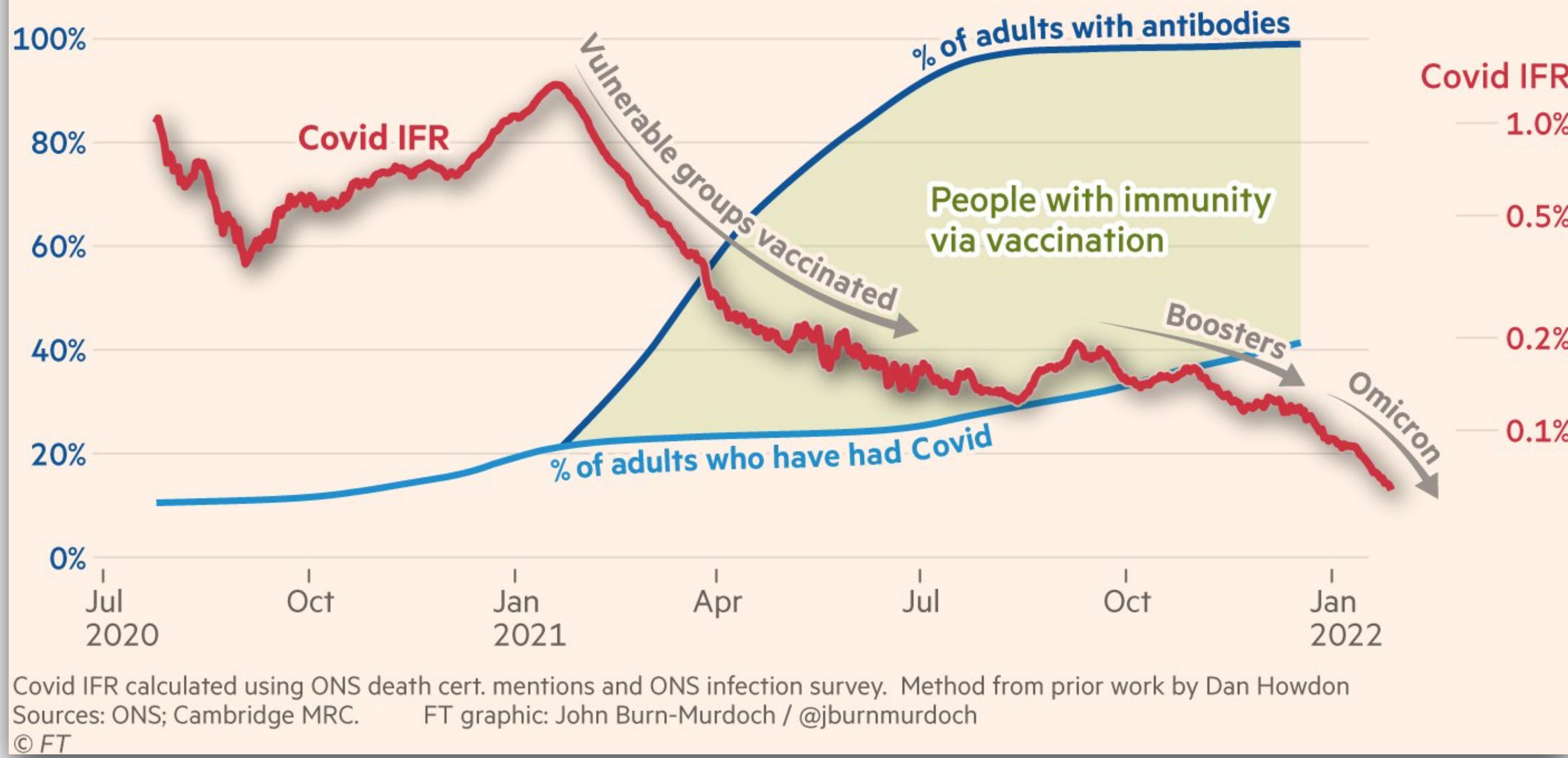
John Tukey

# Data visualisation - Telling a story with data

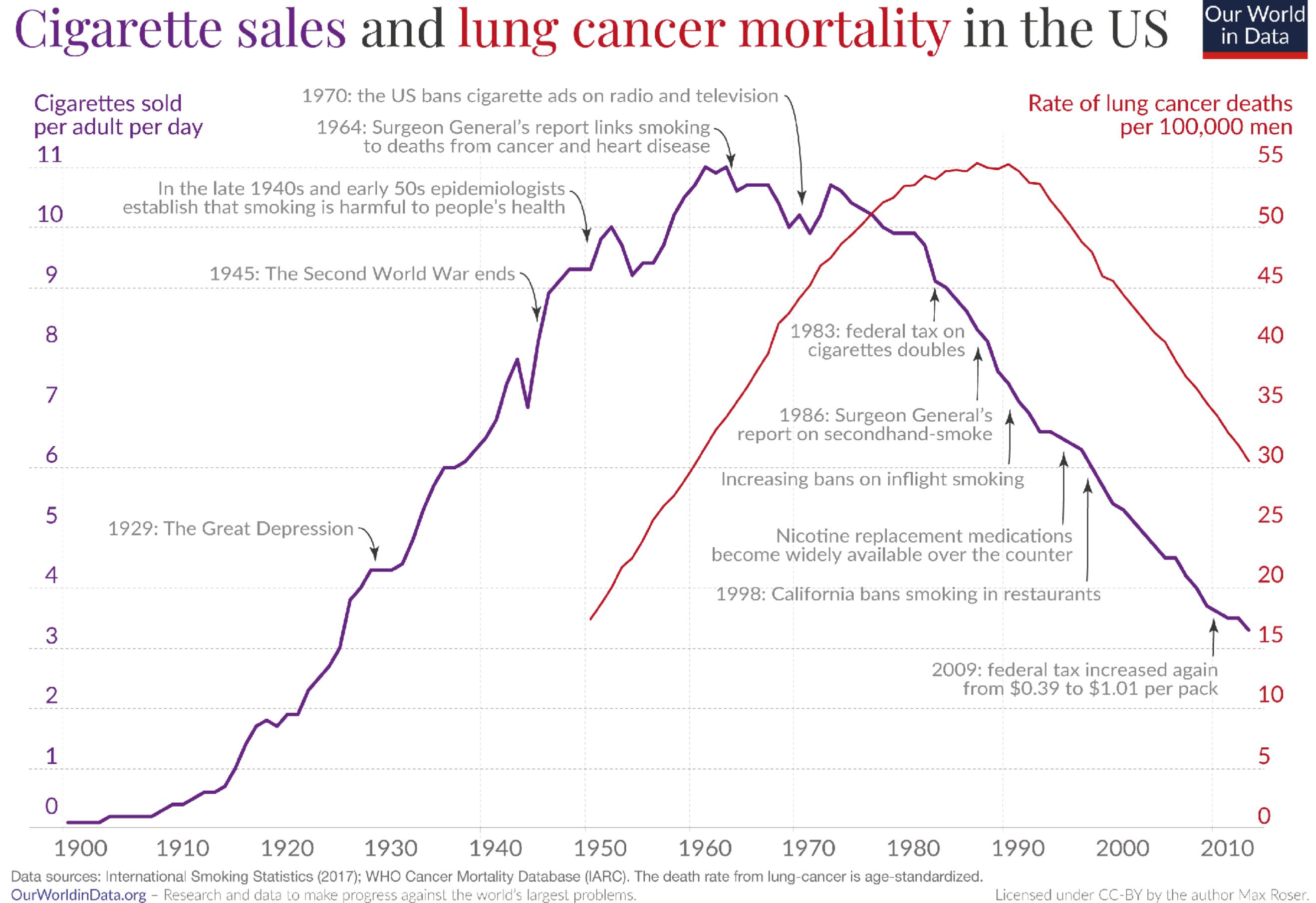
# Data visualisation - Telling a story with data

Covid has grown gradually less lethal over the pandemic, mainly due to immunity, the majority of which has come via vaccines

Evolution of Covid's infection fatality ratio in England, overlaid on levels and sources of immunity



# Data visualisation - Telling a story with data



# Data visualisation - Telling a story with data

## Cigarette sales and lung cancer mortality in the US

Our World  
in Data

Cigarettes sold  
per adult per day

11

1970: the US bans cigarette ads on radio and television

1964: Surgeon General's report links smoking  
to deaths from cancer and heart disease

10

In the late 1940s and early 50s epidemiologists  
establish that smoking is harmful to people's health

9

1945: The

8

7

6

5

1929: The Great Dep

4

3

2

1

0

1900 1910 1920

Rate of lung cancer deaths  
per 100,000 men

55

50

### # of Unique Words Used Within Artist's First 35,000 Lyrics

3,000 words

4,000

5,000

6,000 words

All Just 

Find an Artist

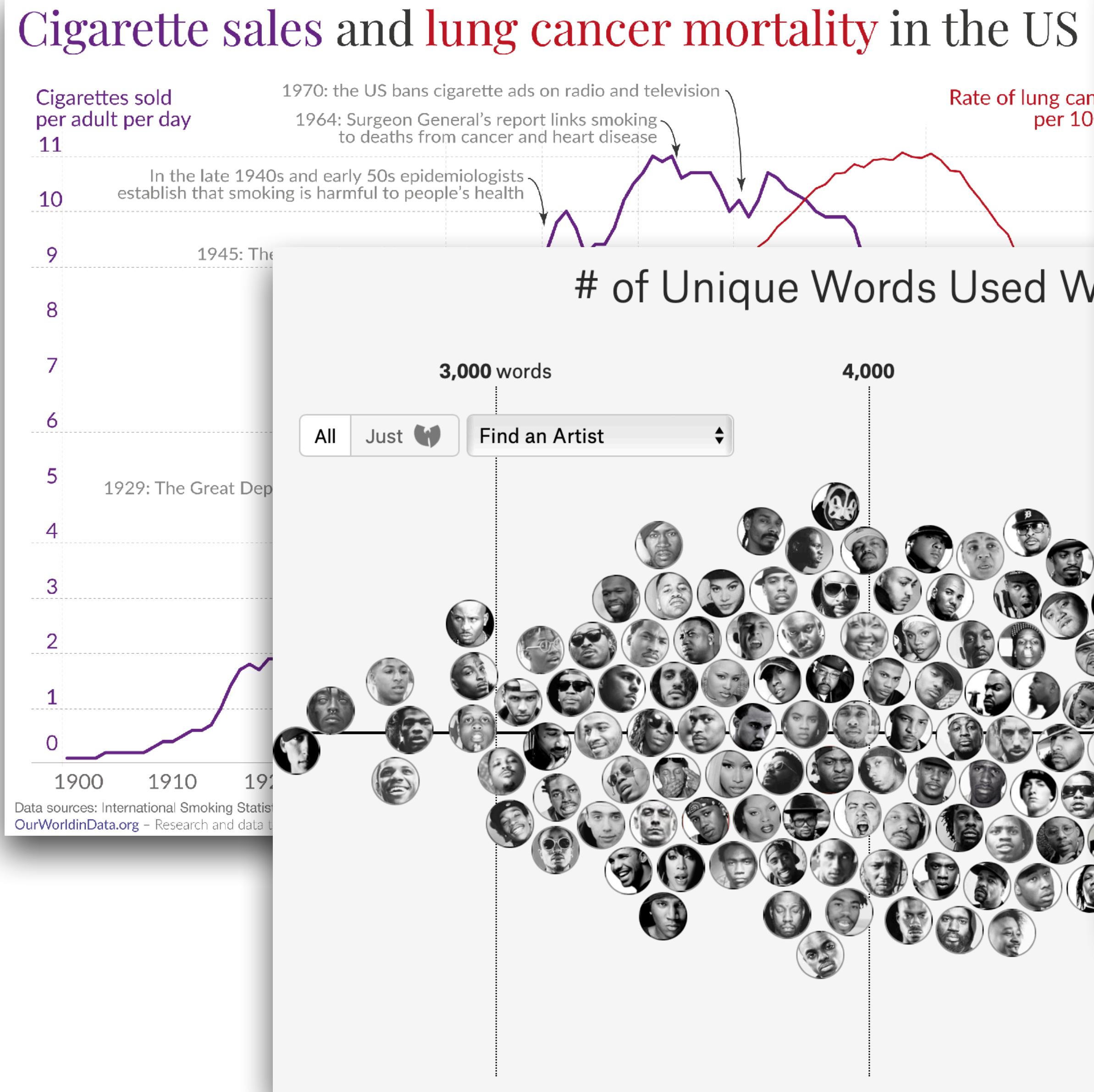
Aesop Rock:  
7,879 unique words used

7,300 words

Data sources: International Smoking Statistics  
OurWorldInData.org – Research and data t

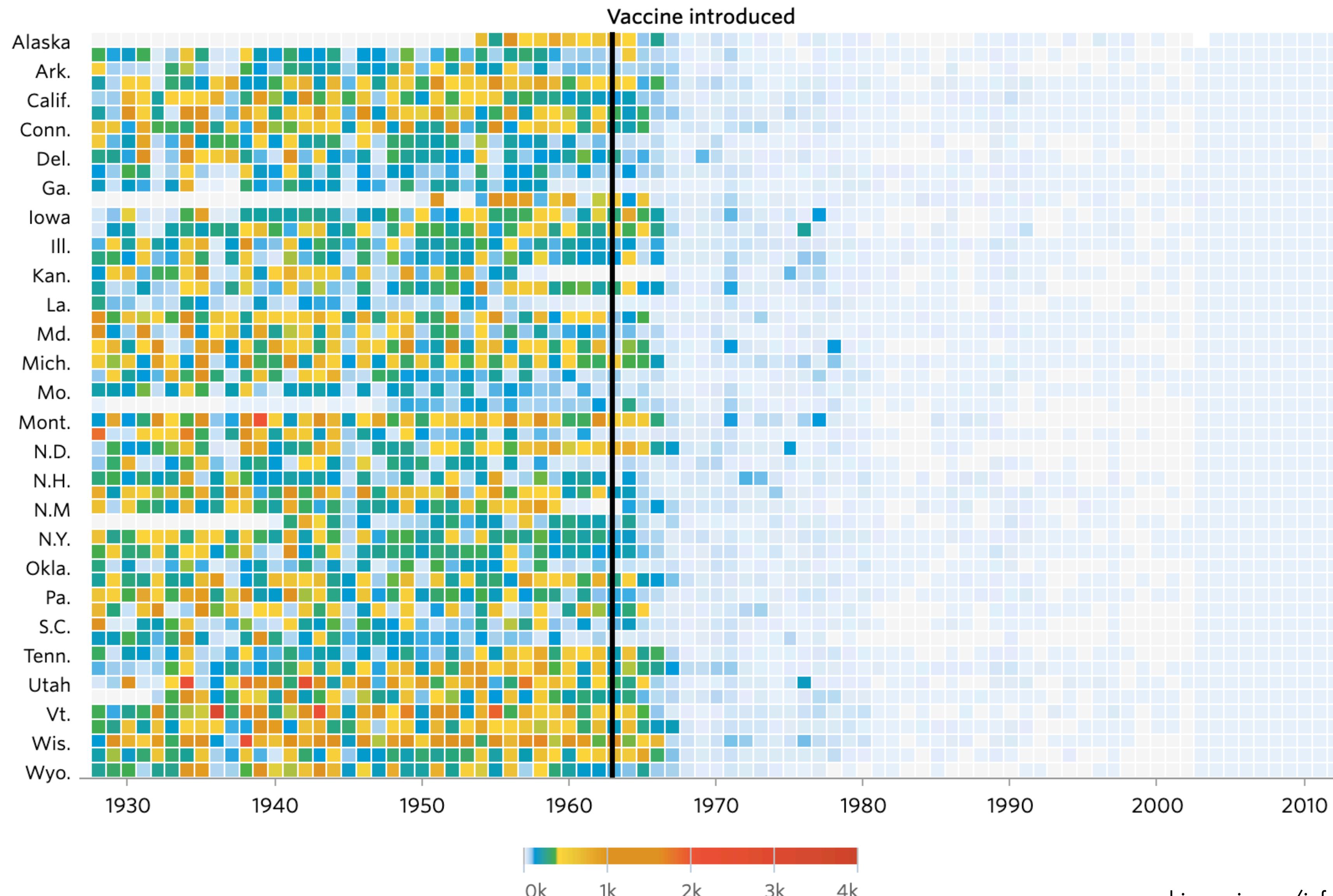
ThePudding

# Data visualisation - Telling a story with data

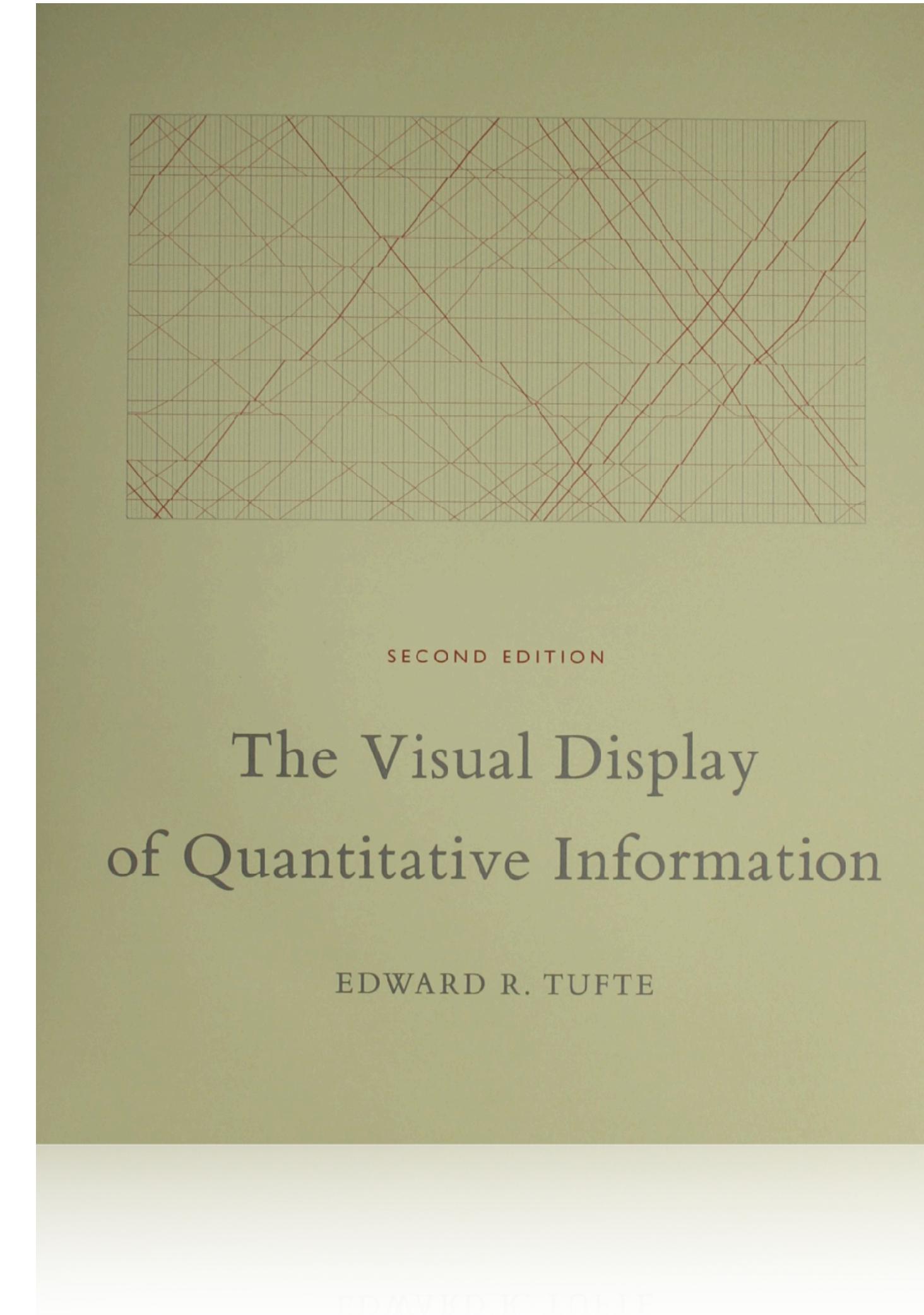


# The impact of vaccines on infectious diseases

## Measles

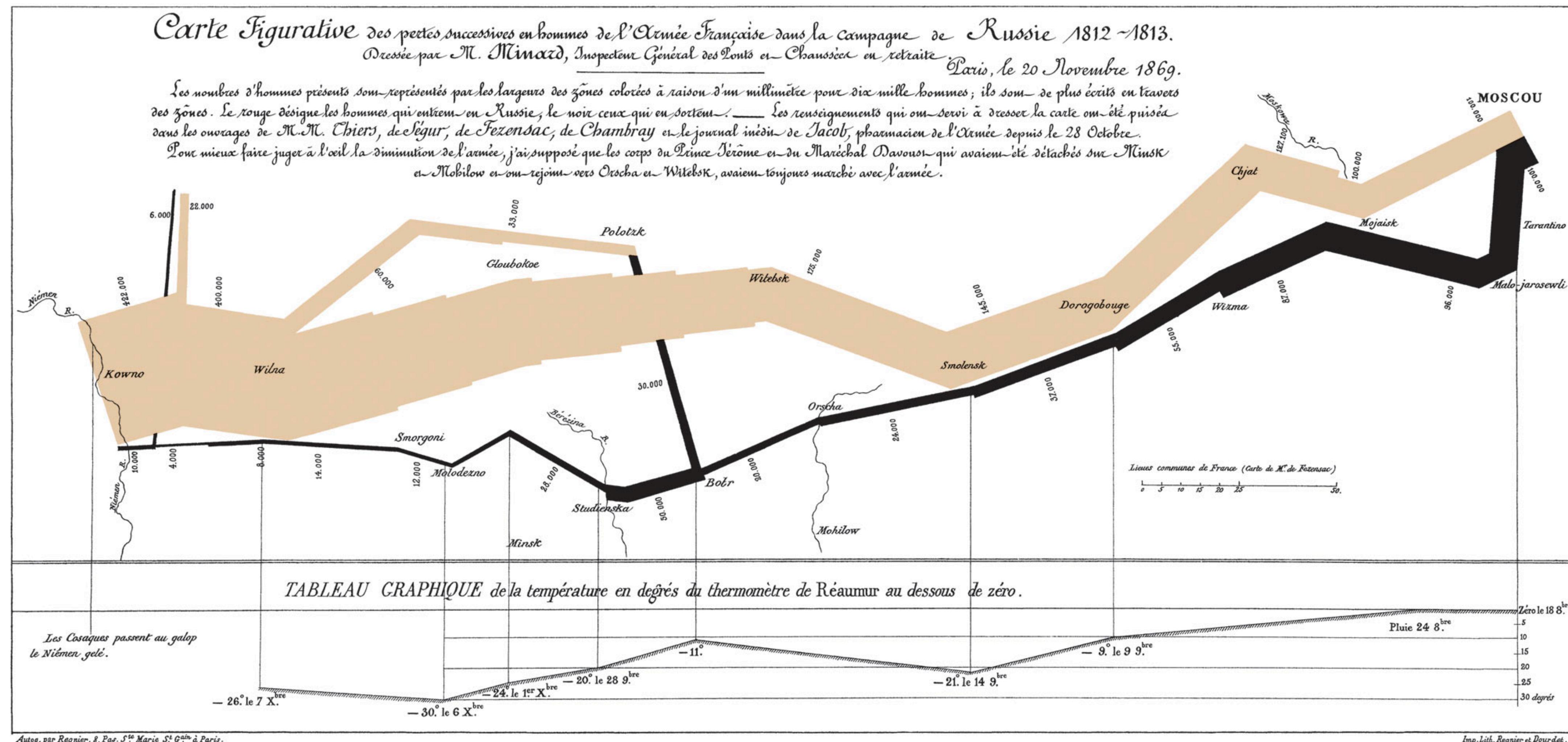


# Pioneers of data visualisation - Edward Tufte ("ET")



[www.edwardtufte.com](http://www.edwardtufte.com)

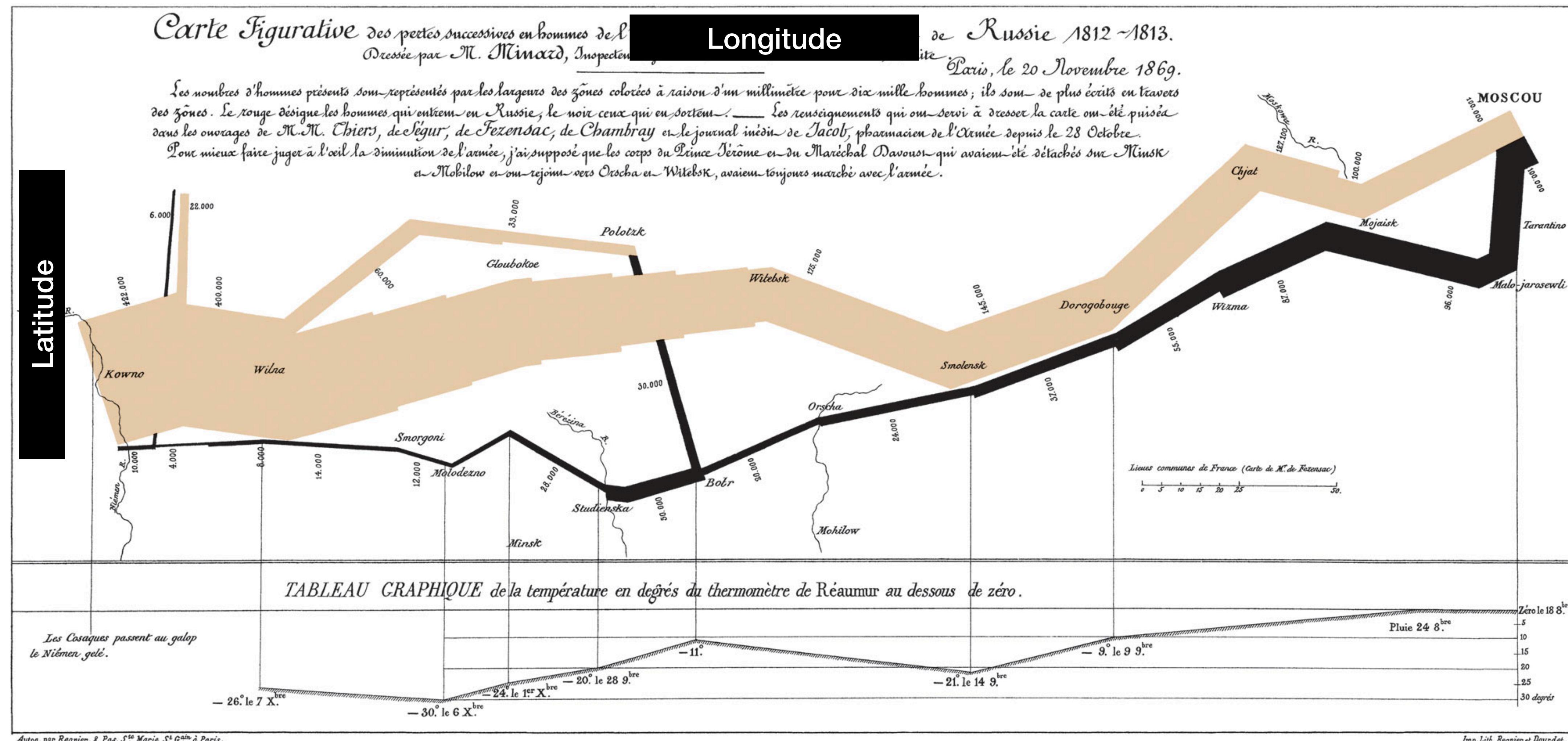
# "The best statistical graph ever drawn"



Charles Joseph Minard

"Carte figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813"

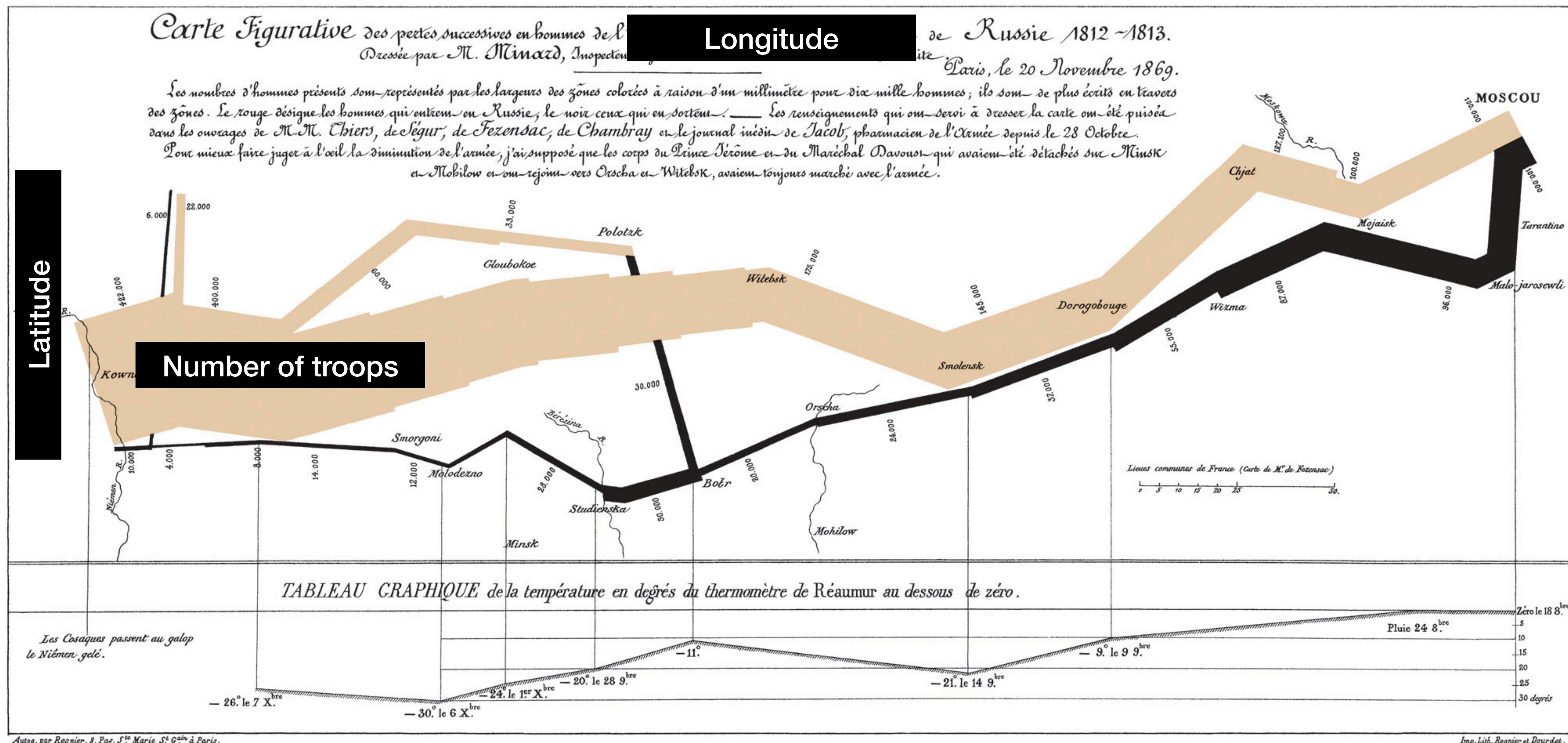
# "The best statistical graph ever drawn"



Charles Joseph Minard

"Carte figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813"

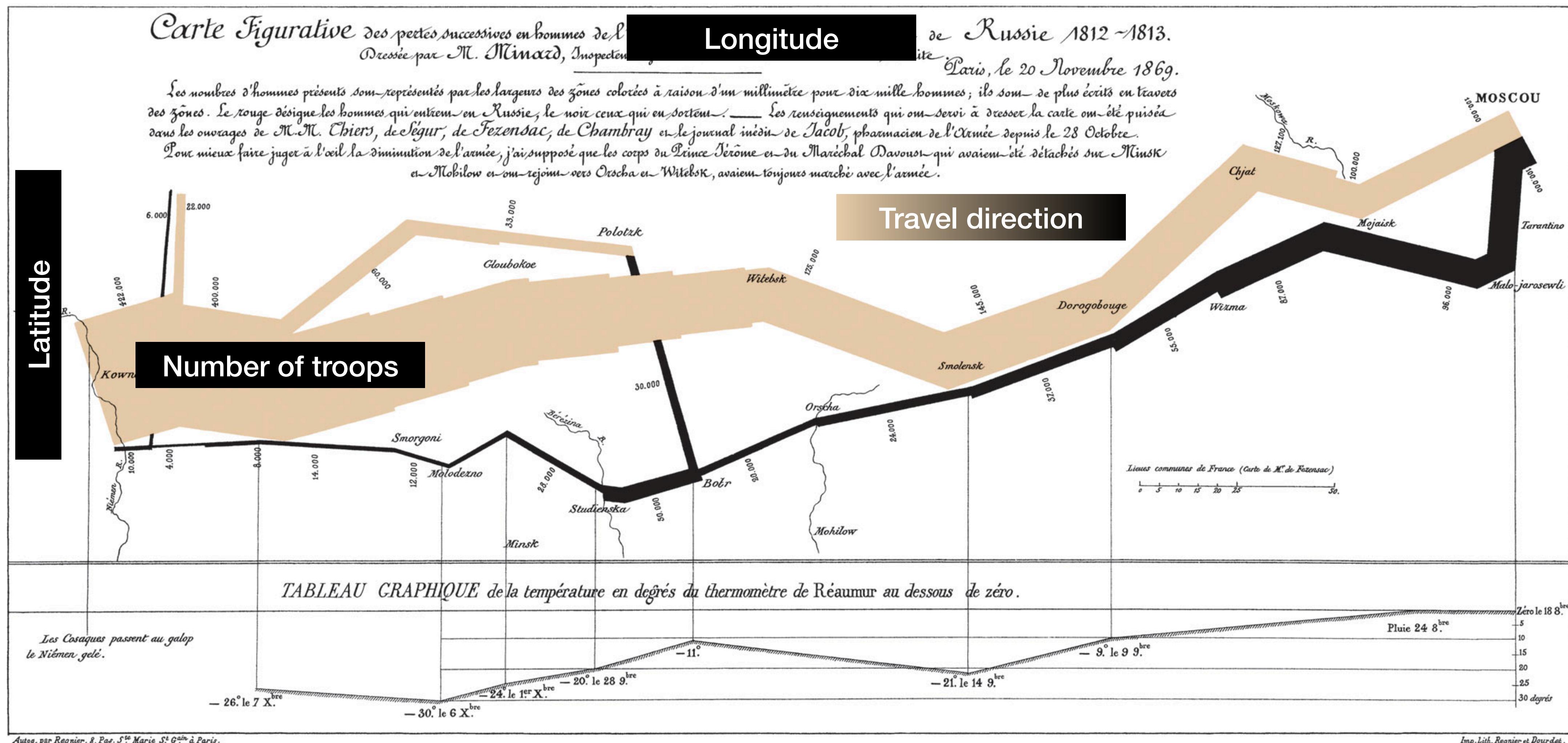
# "The best statistical graph ever drawn"



Charles Joseph Minard

"Carte figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813"

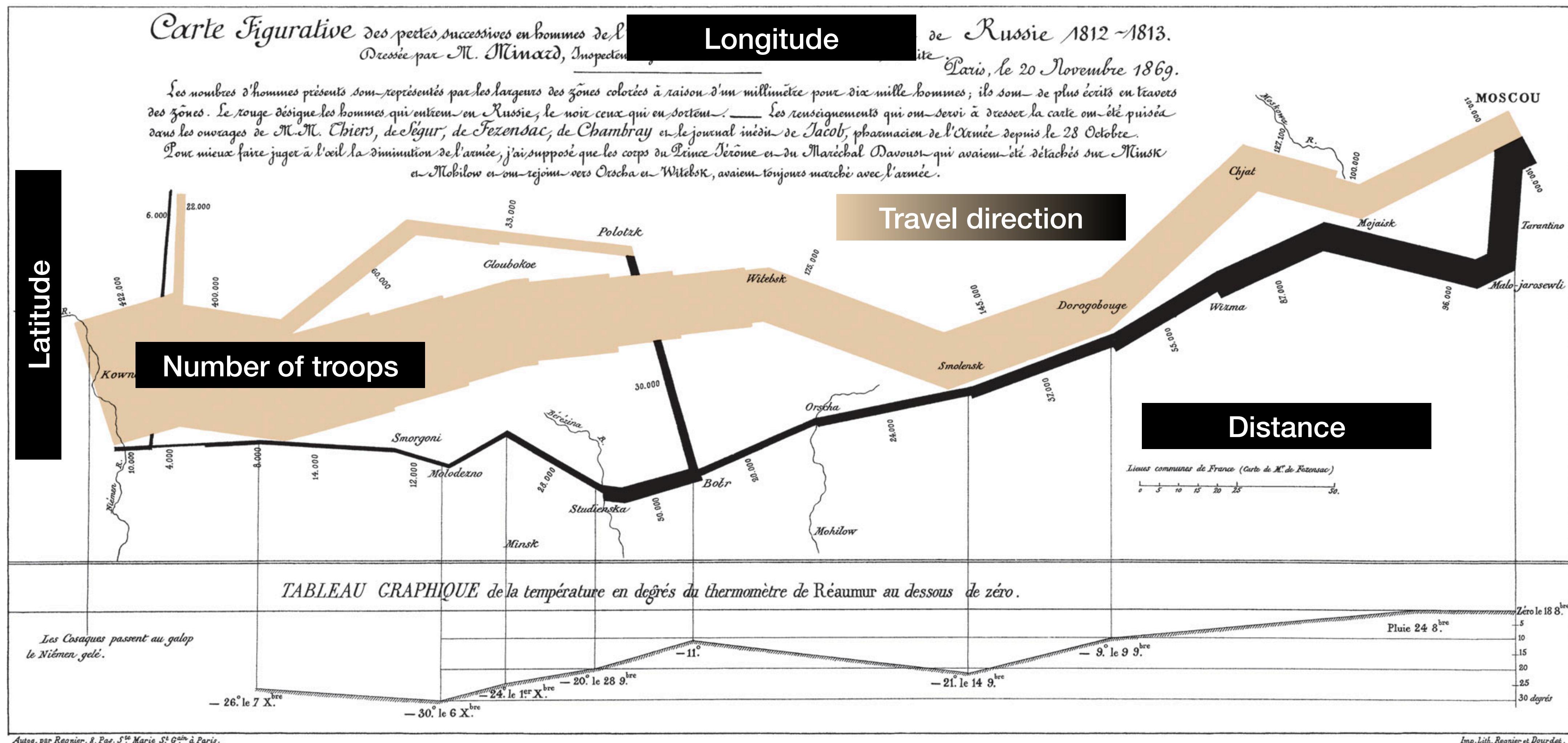
# "The best statistical graph ever drawn"



Charles Joseph Minard

"Carte figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812–1813"

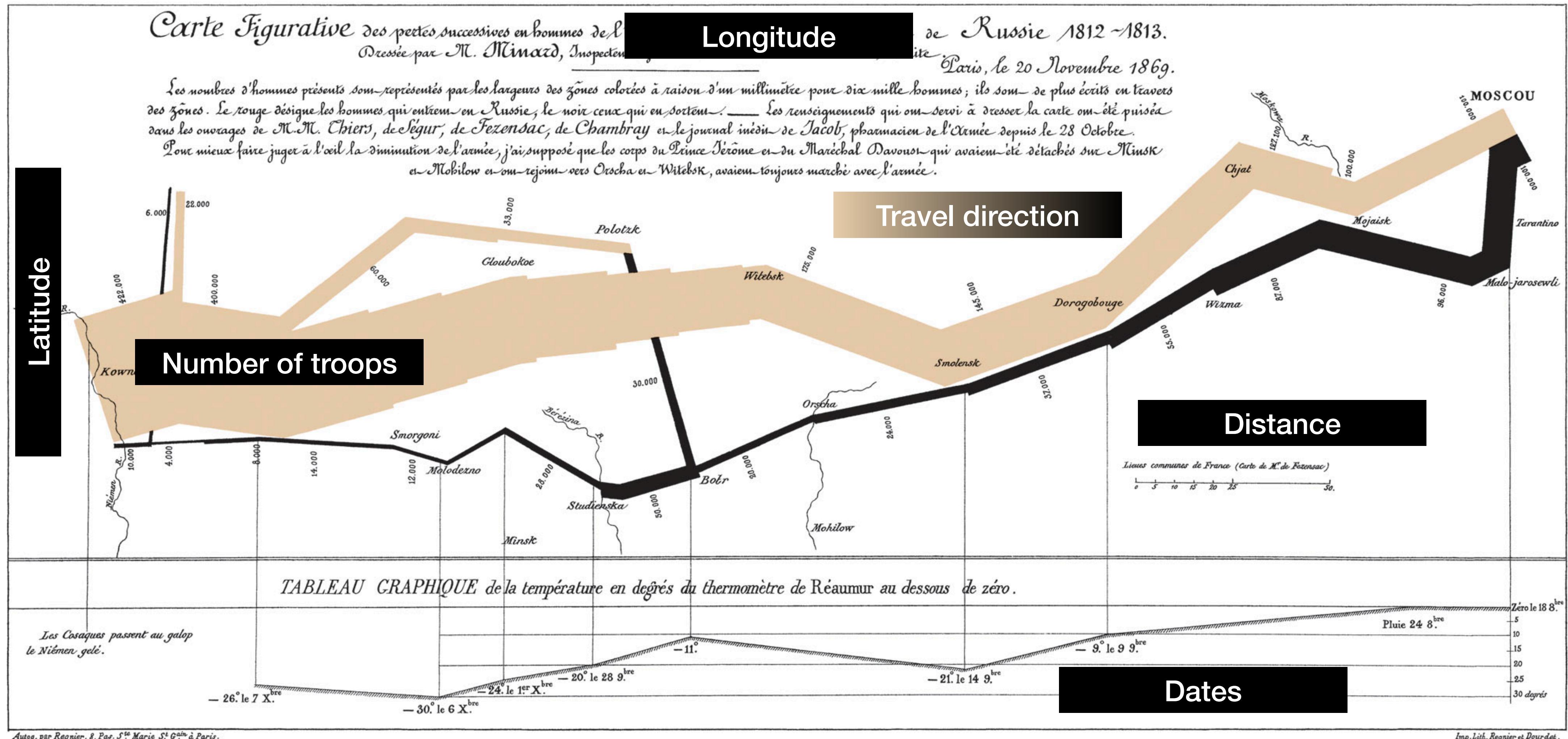
# "The best statistical graph ever drawn"



Charles Joseph Minard

"Carte figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812–1813"

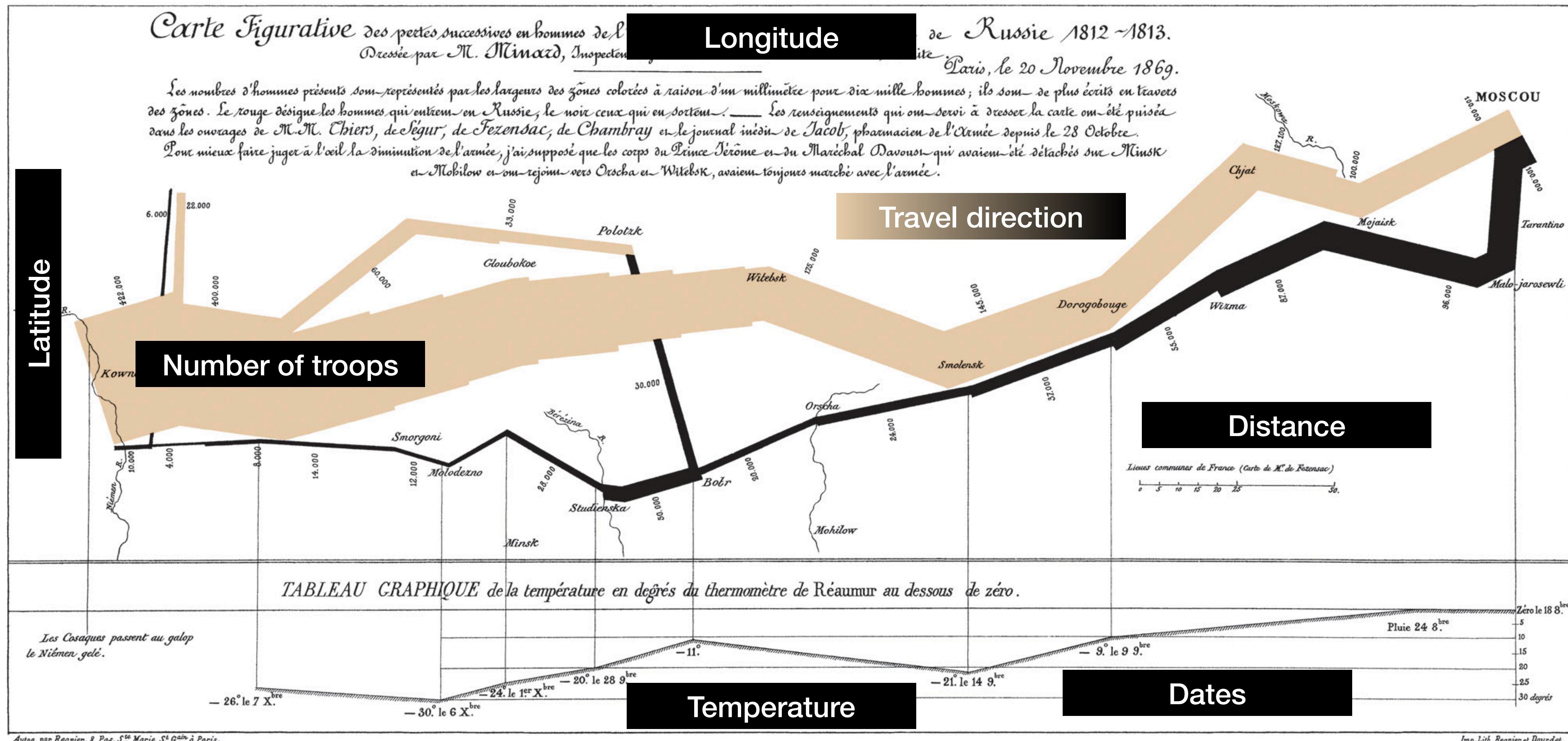
# "The best statistical graph ever drawn"



Charles Joseph Minard

"Carte figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812–1813"

# "The best statistical graph ever drawn"

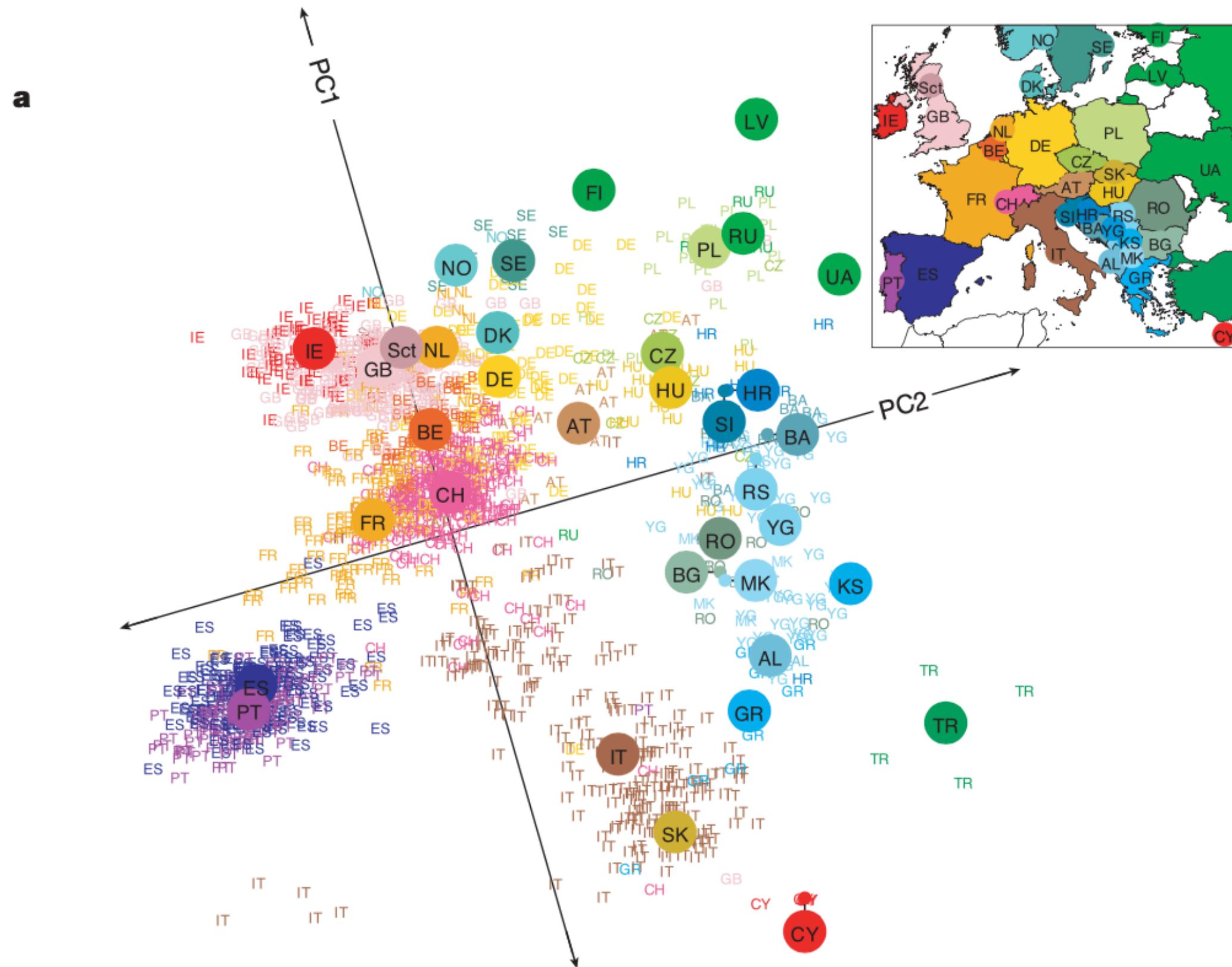


Charles Joseph Minard

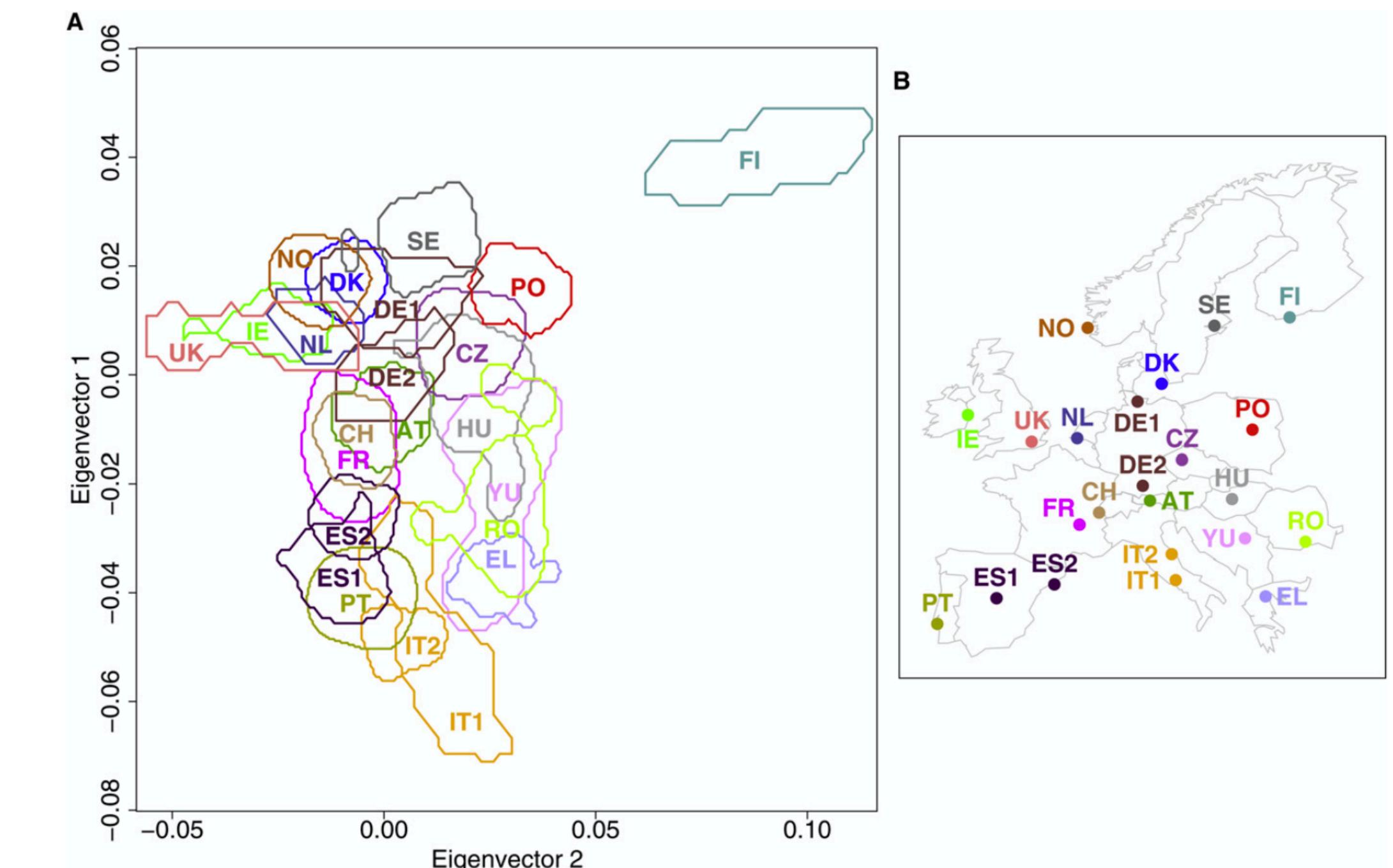
"Carte figurative des pertes successives en hommes de l'Armée Française dans la campagne de Russie 1812-1813"

# A tale of two visualisations

Novembre et al 2008



Lao et al 2008



# Which visualisation do you prefer?

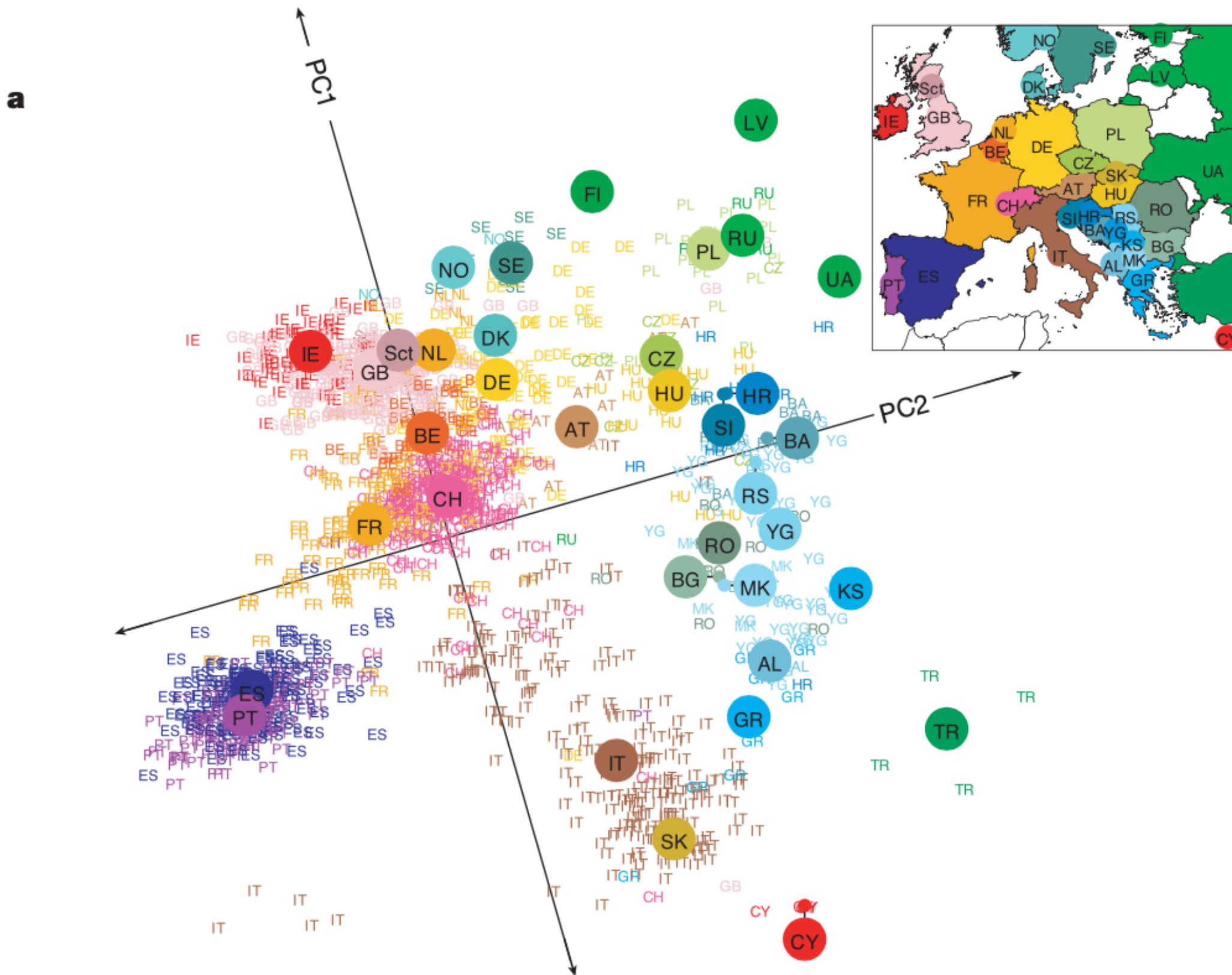
# A tale of two visualisations

Novembre et al, 31/8/2008

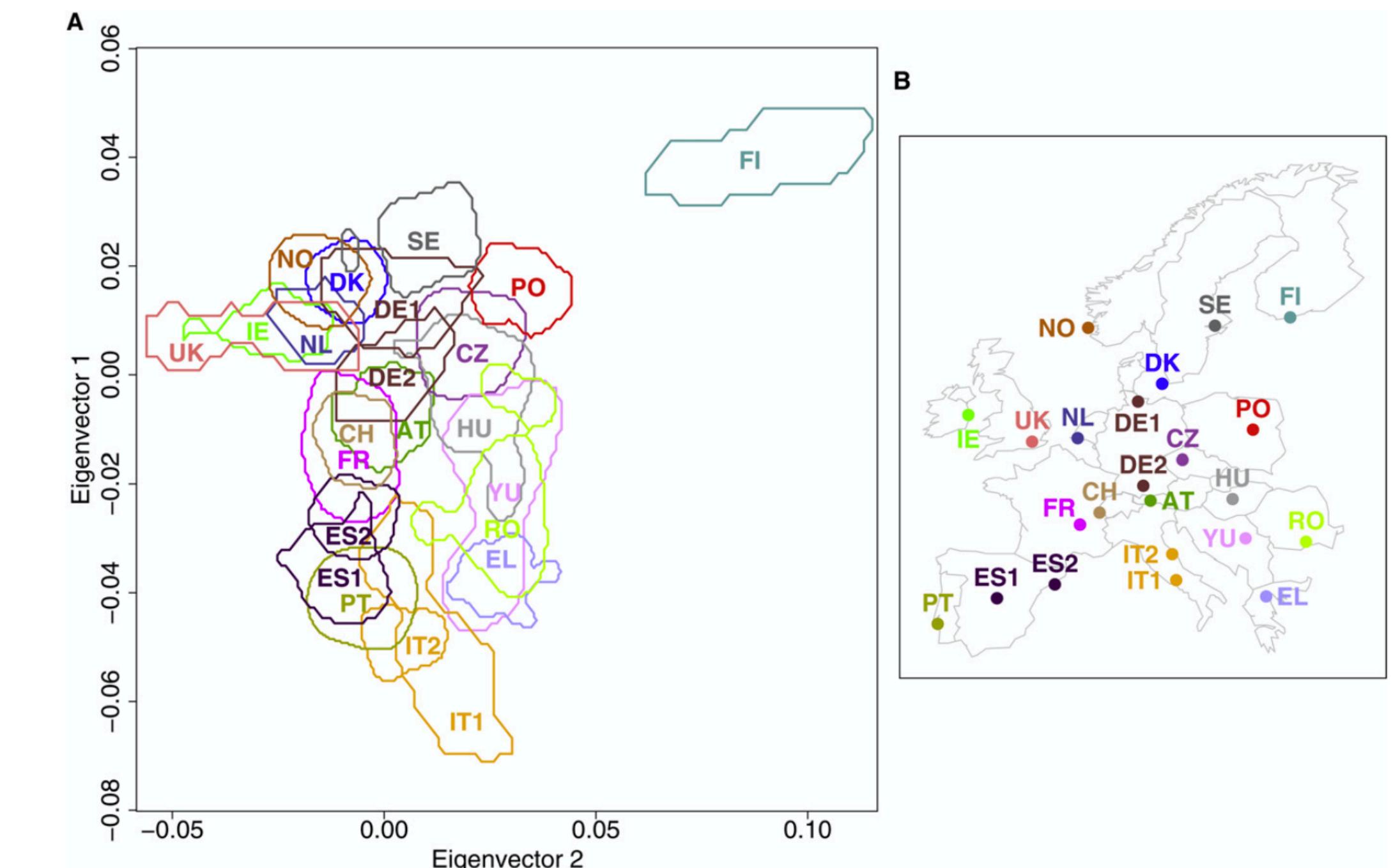
# “Genes mirror geography in Europe”

Lao et al, 26/8/2008

# "Correlation between Genetic and Geographic Structure in Europe"

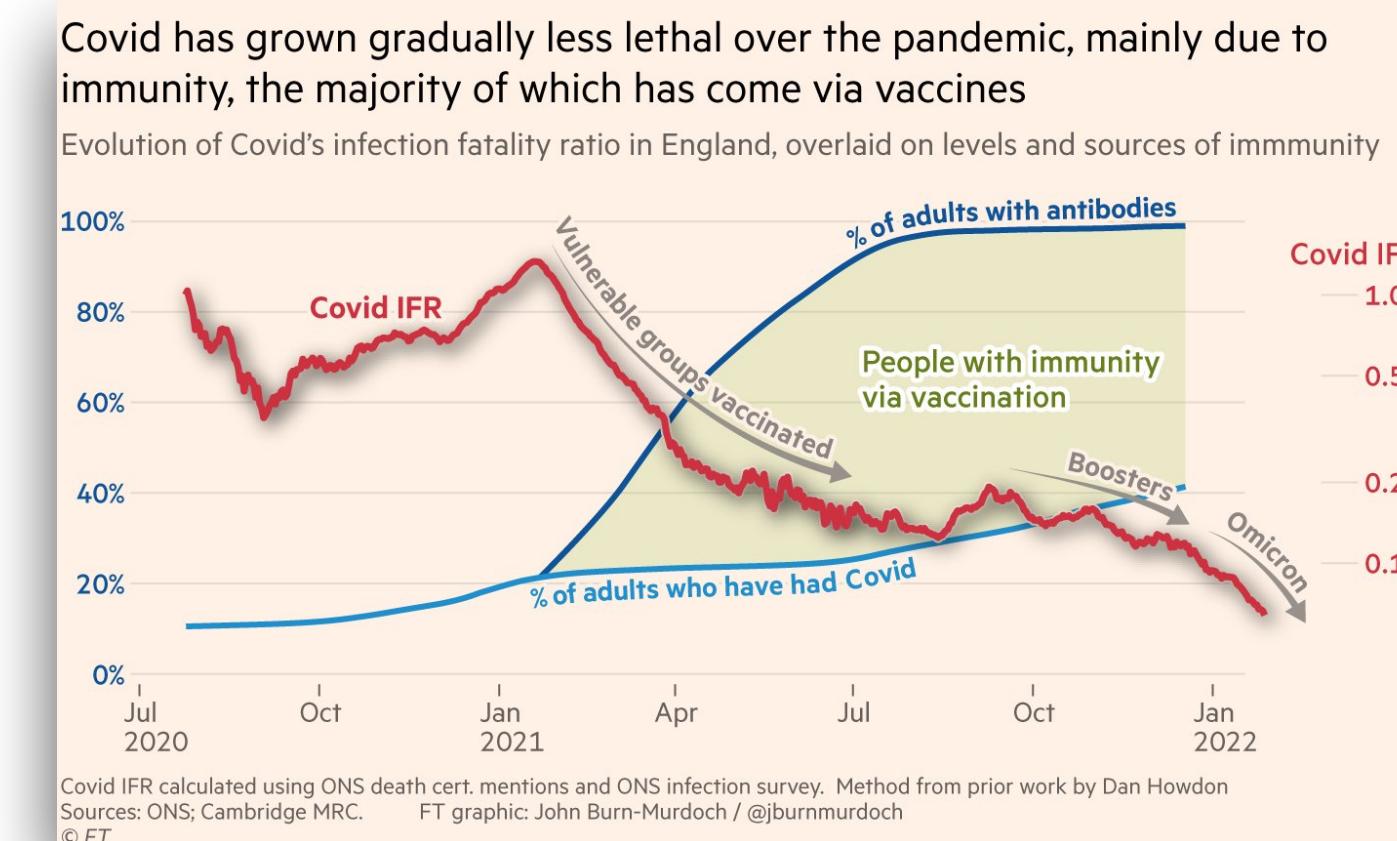
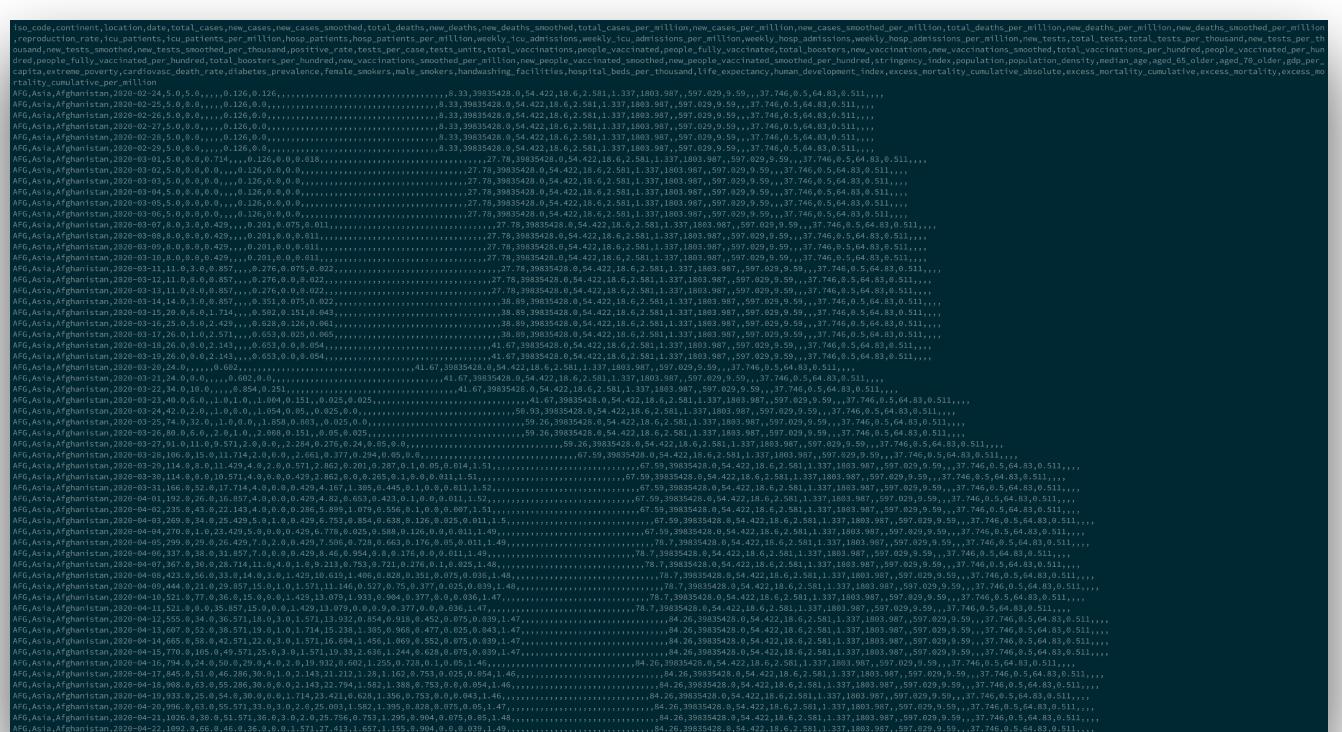


1,883 citations (Nature)

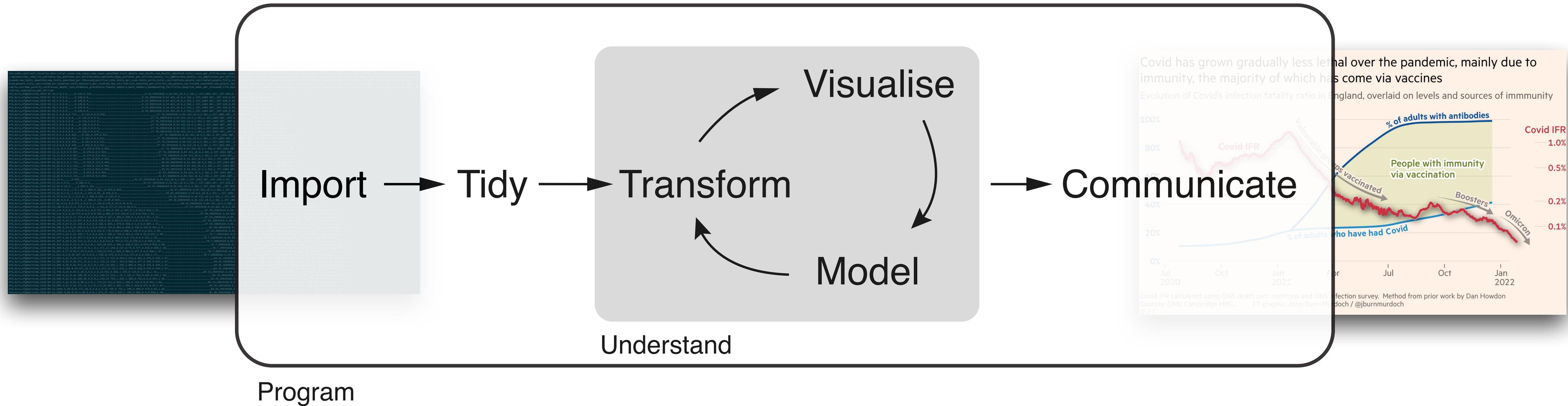


630 citations (Current Biology)

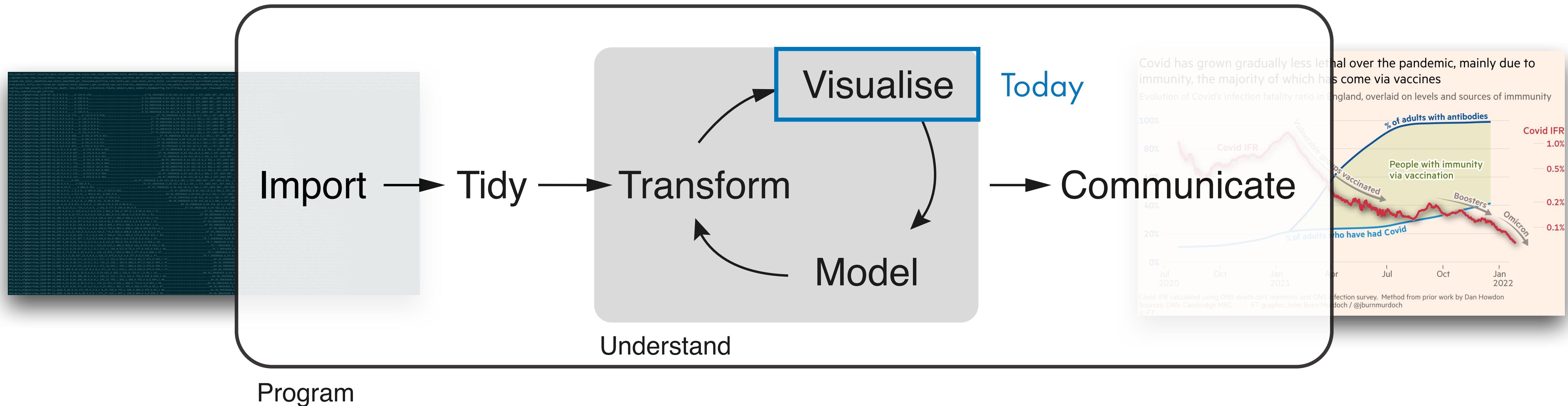
# A data science project workflow



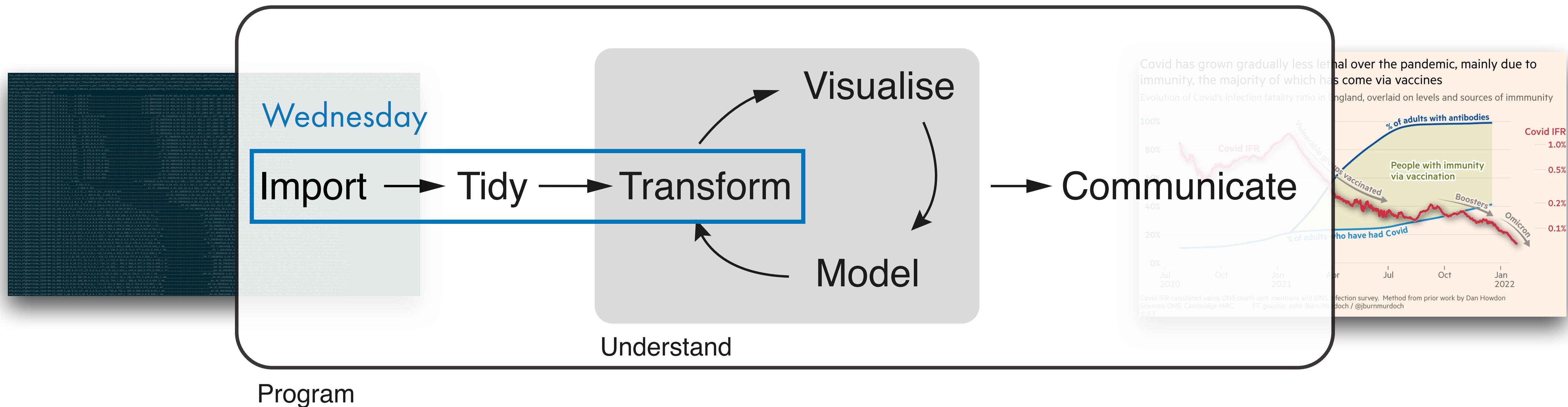
# A data science project workflow



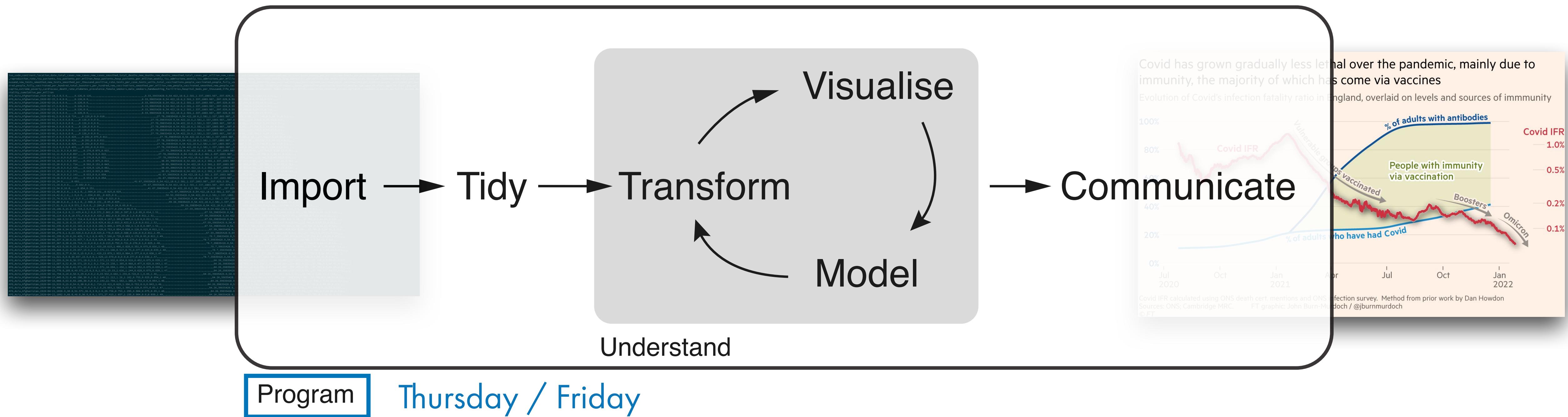
# A data science project workflow



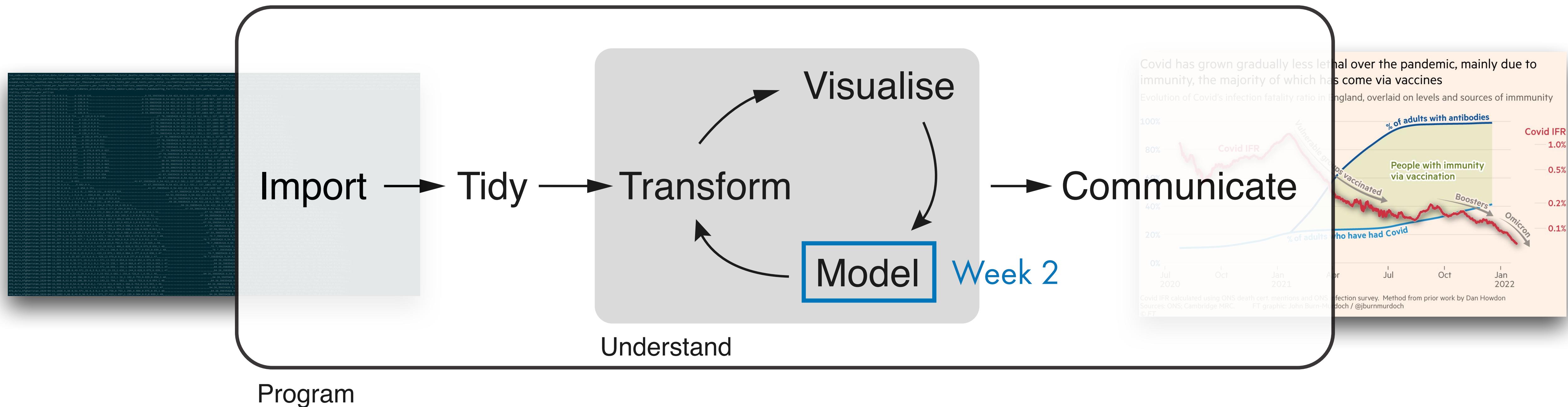
# A data science project workflow



# A data science project workflow



# A data science project workflow



# The tidyverse



Tidyverse

Packages   Blog   Learn   Help   Contribute

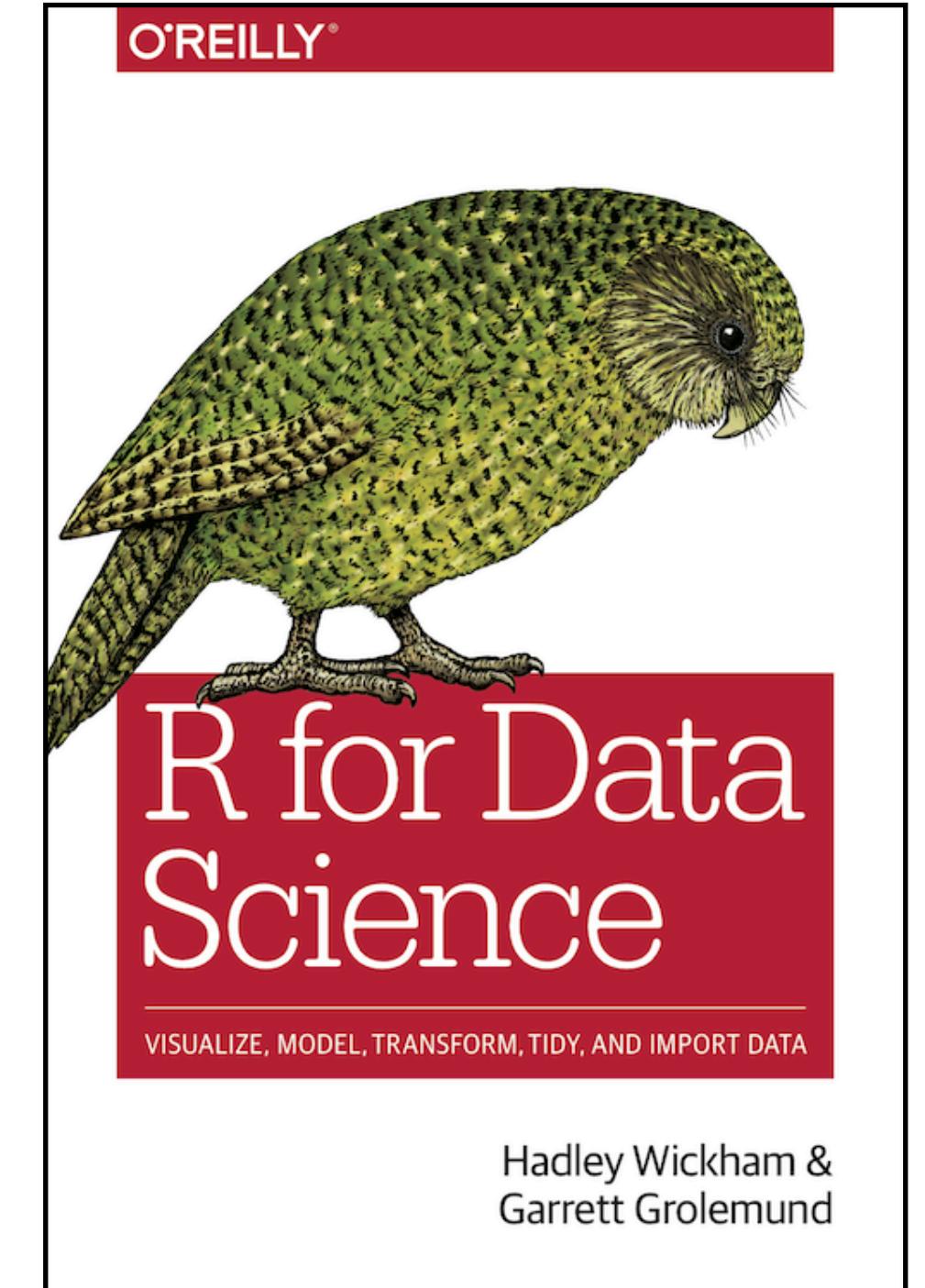
R packages for data science

The tidyverse is an opinionated **collection of R packages** designed for data science. All packages share an underlying design philosophy, grammar, and data structures.

Install the complete tidyverse with:

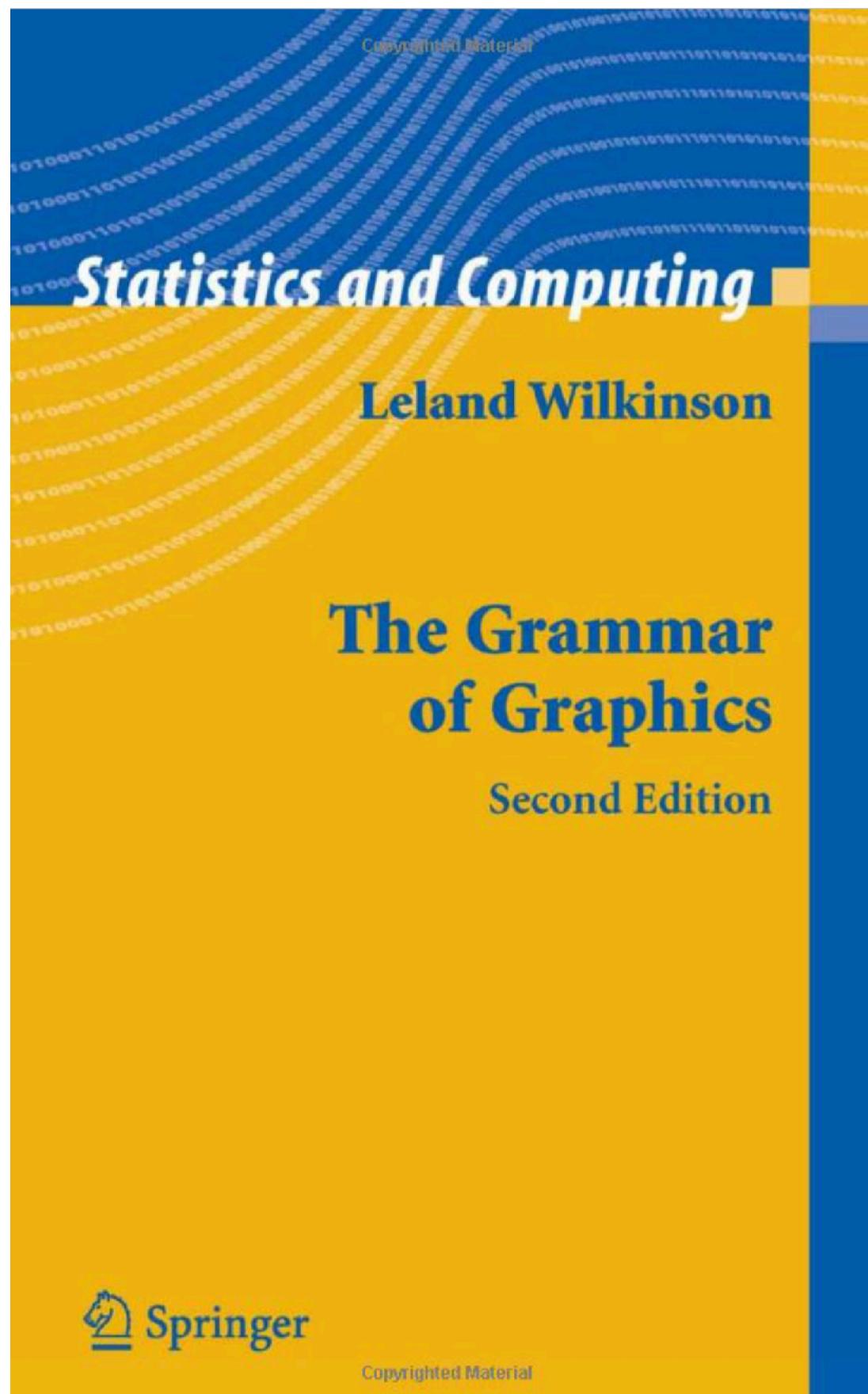
```
install.packages("tidyverse")
```

[www.tidyverse.org](http://www.tidyverse.org)

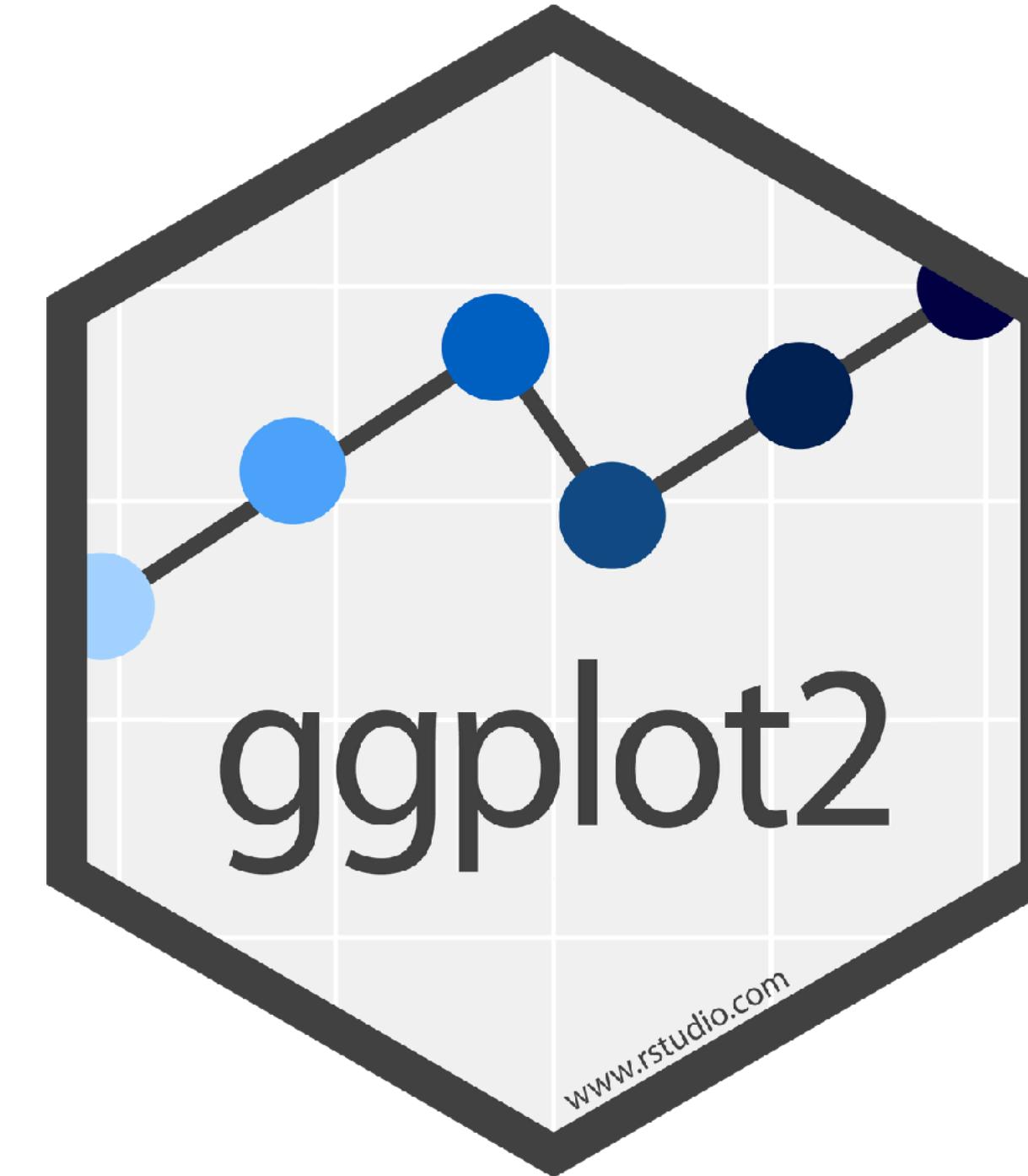


[r4ds.had.co.nz](http://r4ds.had.co.nz)

# A layered grammar of graphics



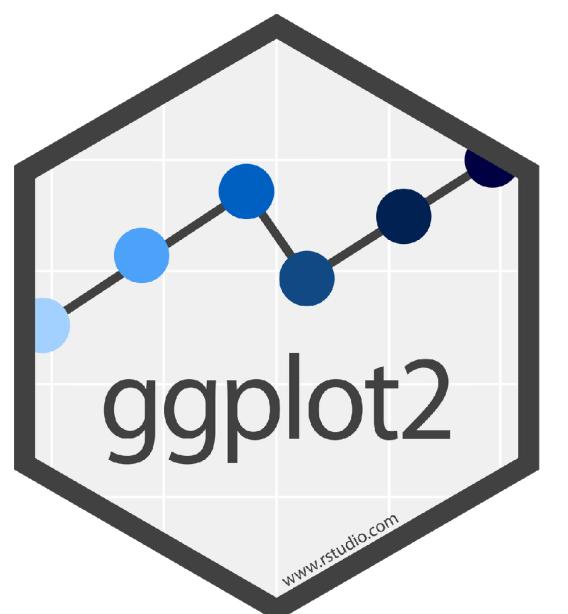
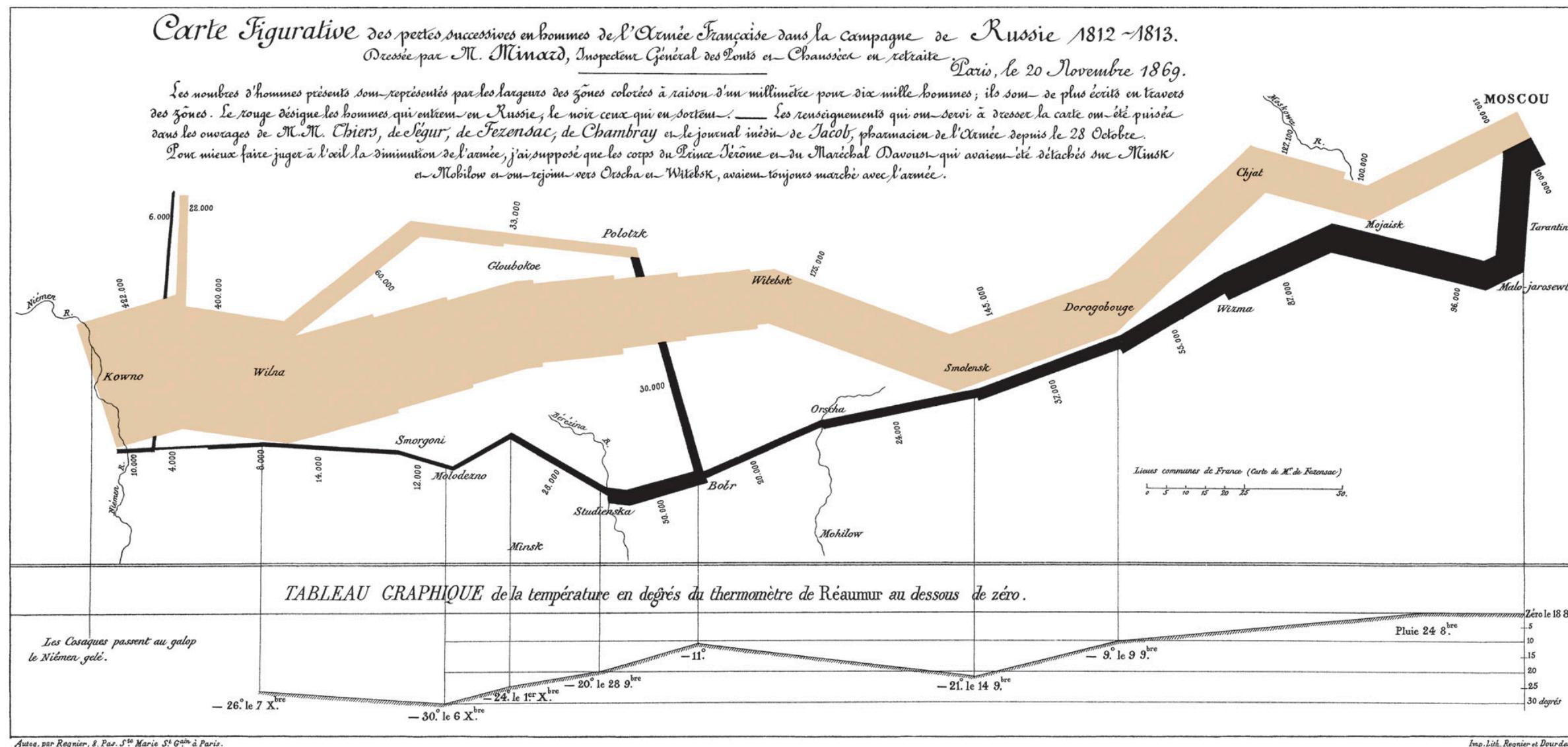
A description of the fundamental features  
that underlie all statistical graphics



A toolkit implementing the grammar of  
graphics in R

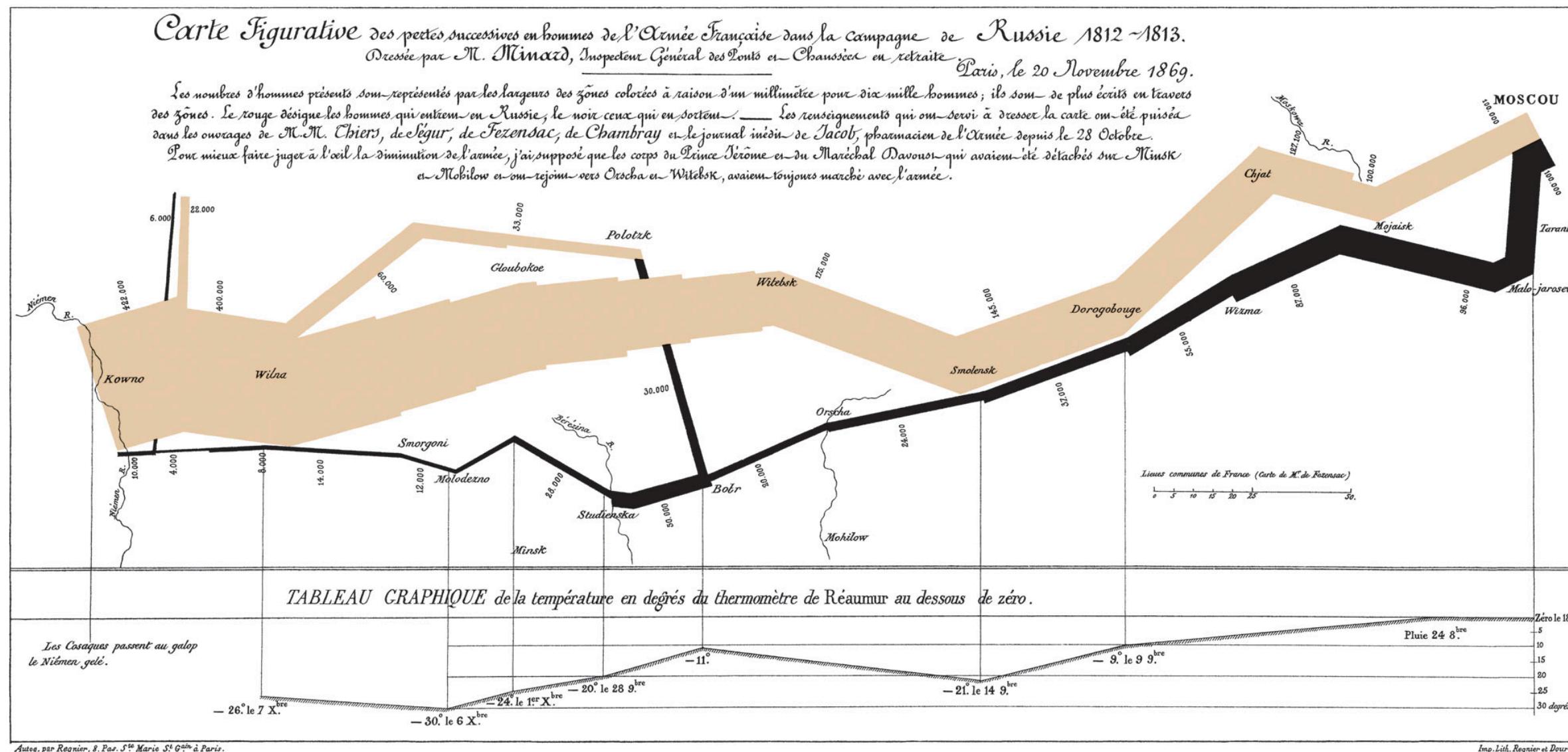
# A layered grammar of graphics

Original

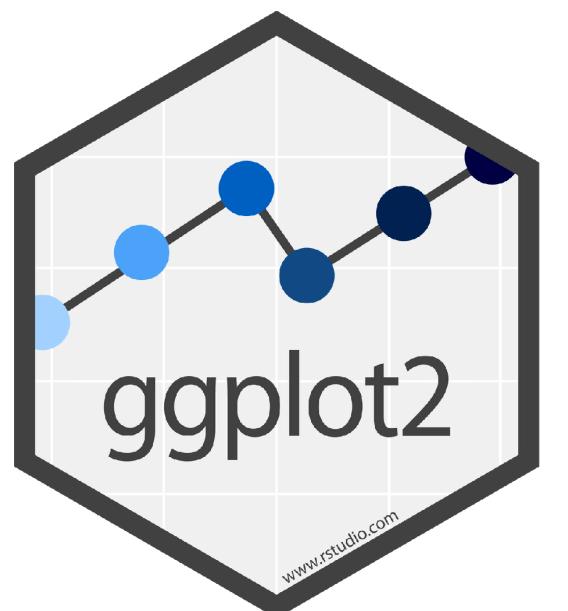
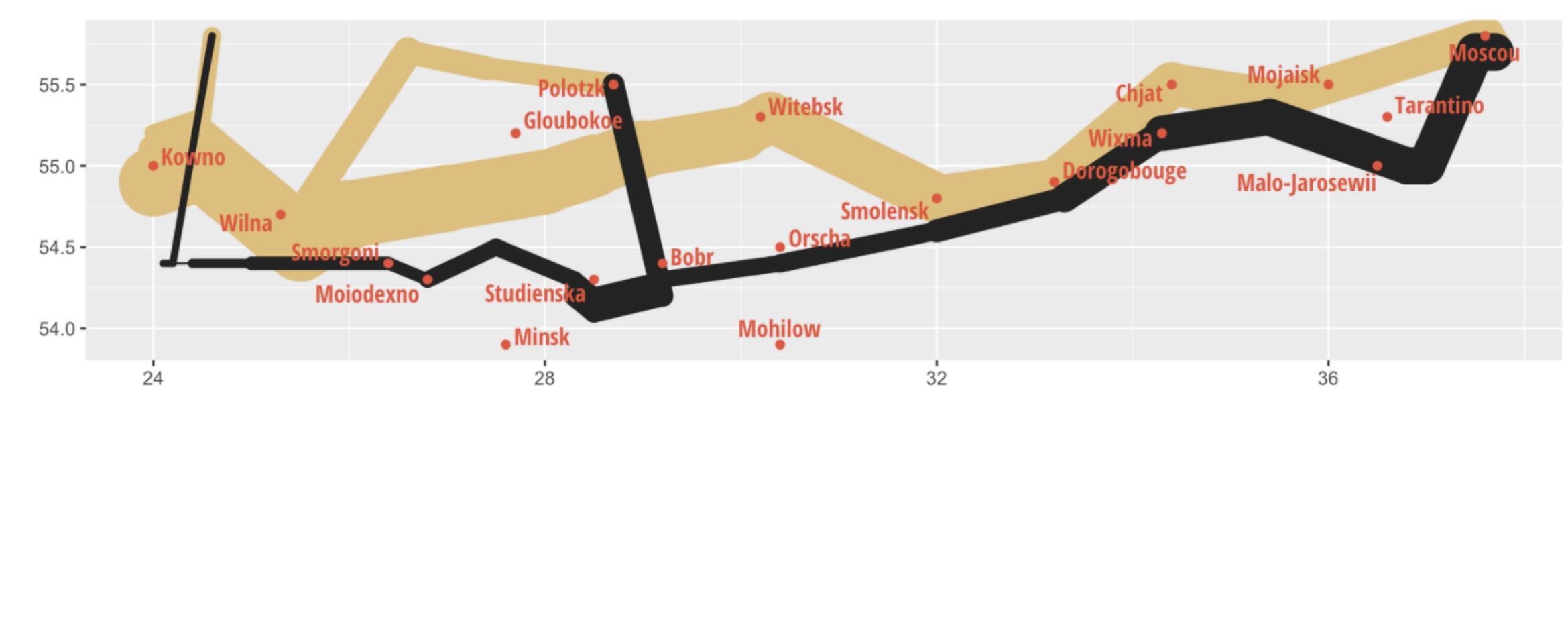


# A layered grammar of graphics

Original

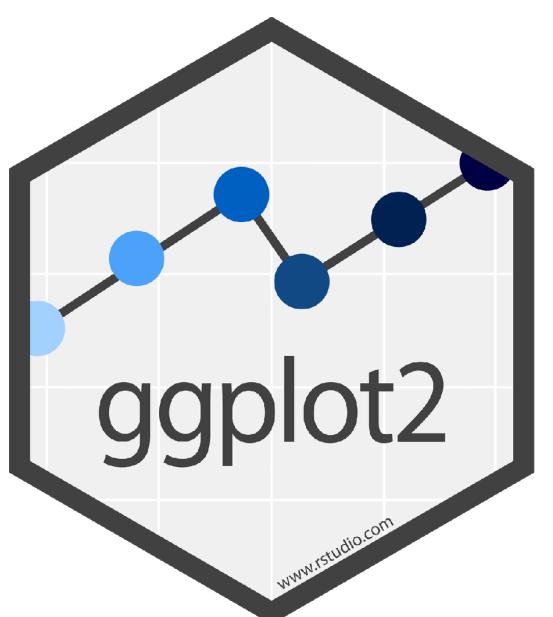
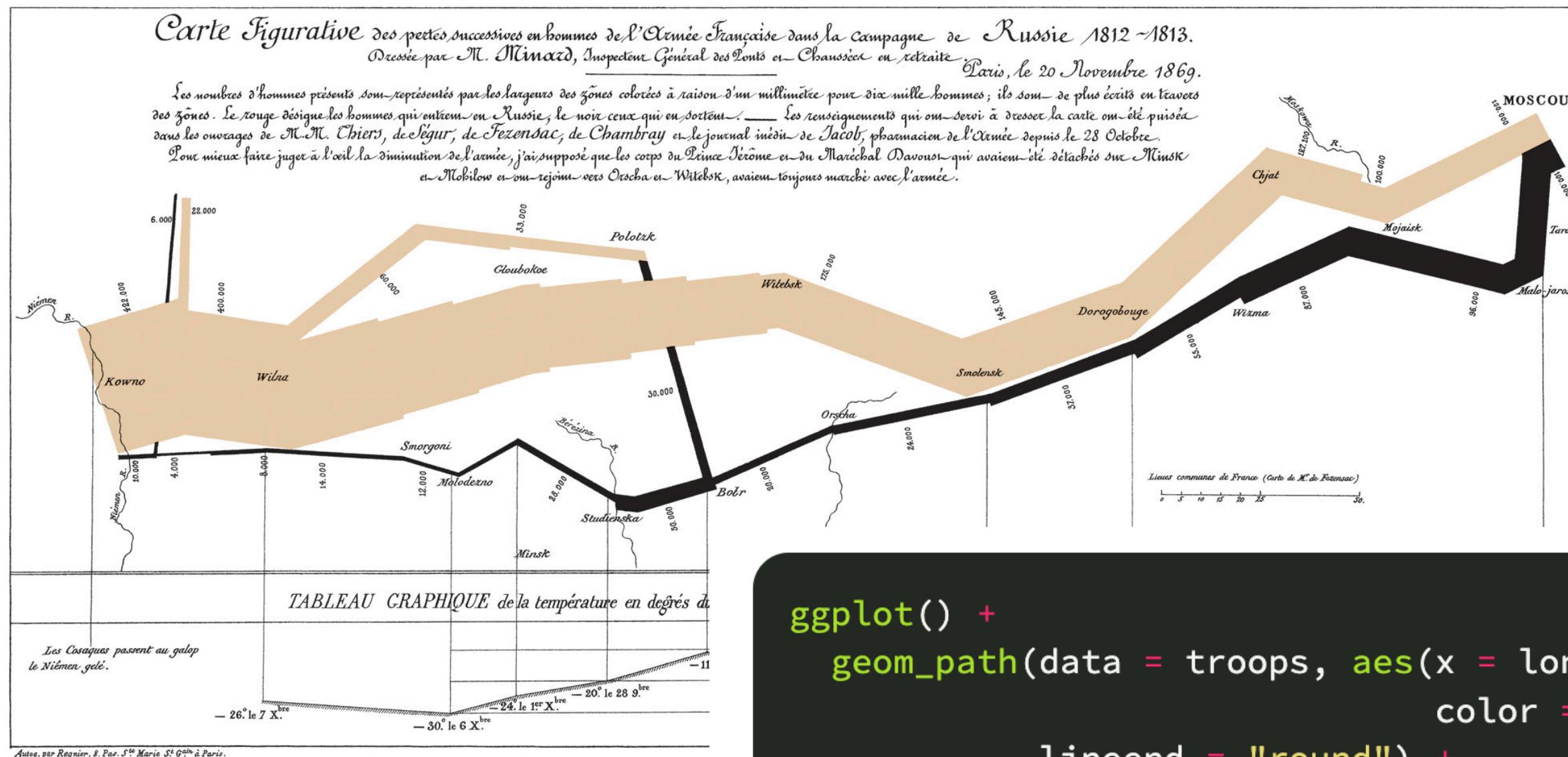


ggplot2 version

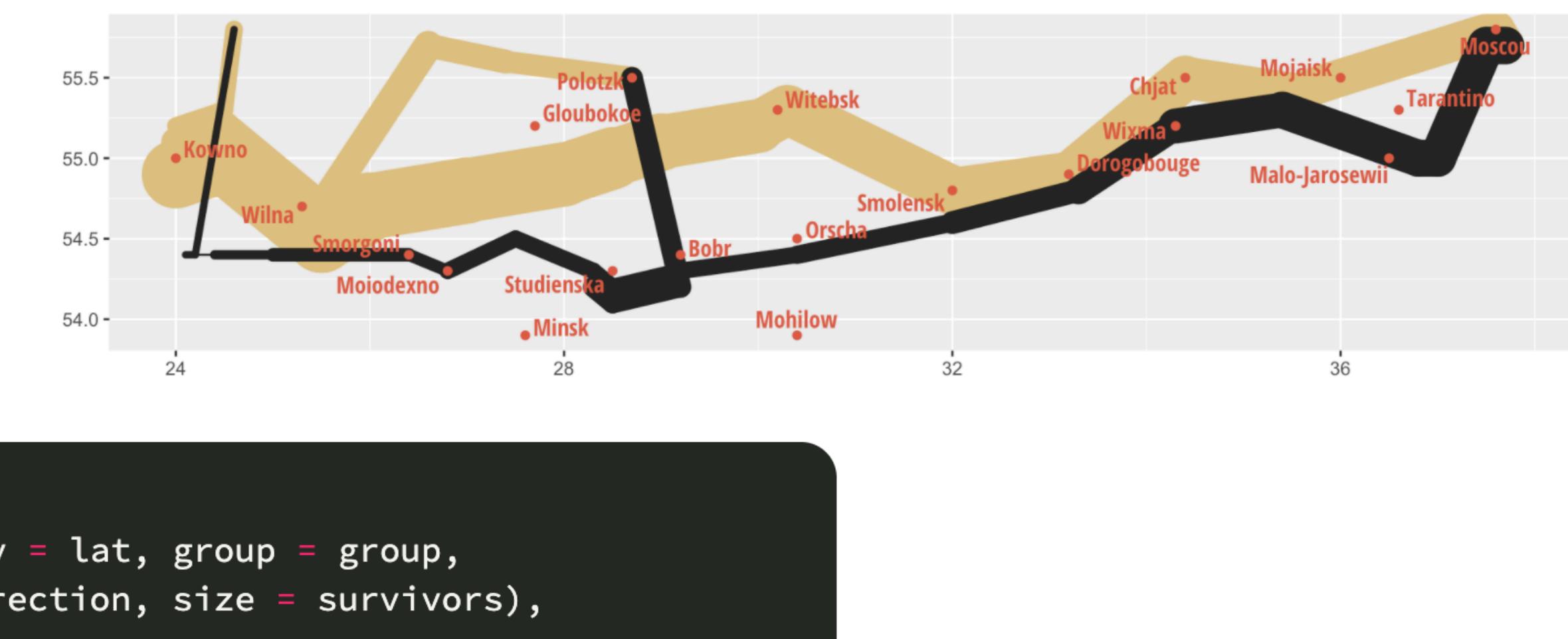


# A layered grammar of graphics

Original



ggplot2 version



```
ggplot() +  
  geom_path(data = troops, aes(x = long, y = lat, group = group,  
                               color = direction, size = survivors),  
            lineend = "round") +  
  geom_point(data = cities, aes(x = long, y = lat),  
             color = "#DC5B44") +  
  geom_text_repel(data = cities, aes(x = long, y = lat, label = city),  
                 color = "#DC5B44", family = "Open Sans Condensed Bold") +  
  scale_size(range = c(0.5, 15)) +  
  scale_colour_manual(values = c("#DFC17E", "#252523")) +  
  labs(x = NULL, y = NULL) +  
  guides(color = FALSE, size = FALSE)
```

# A layered grammar of graphics

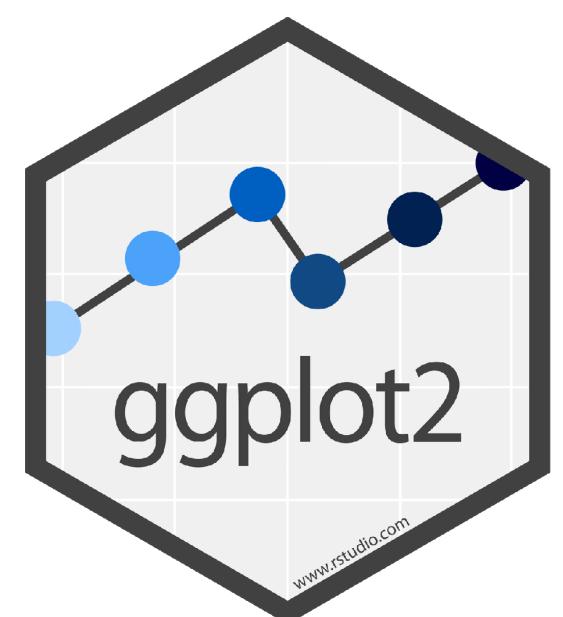
data

MAPPED  
TO

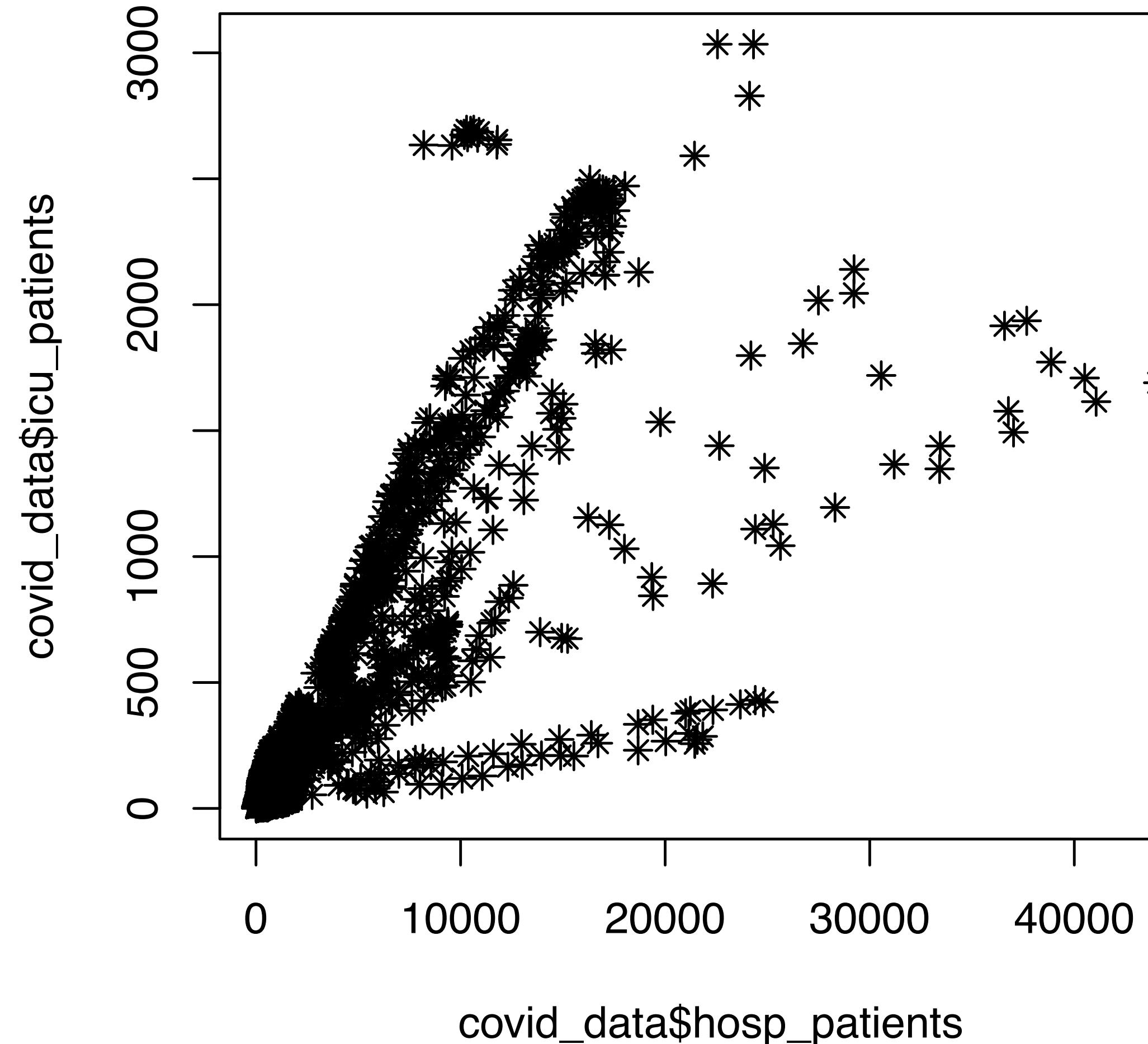
aesthetic  
properties

OF

geometric  
objects



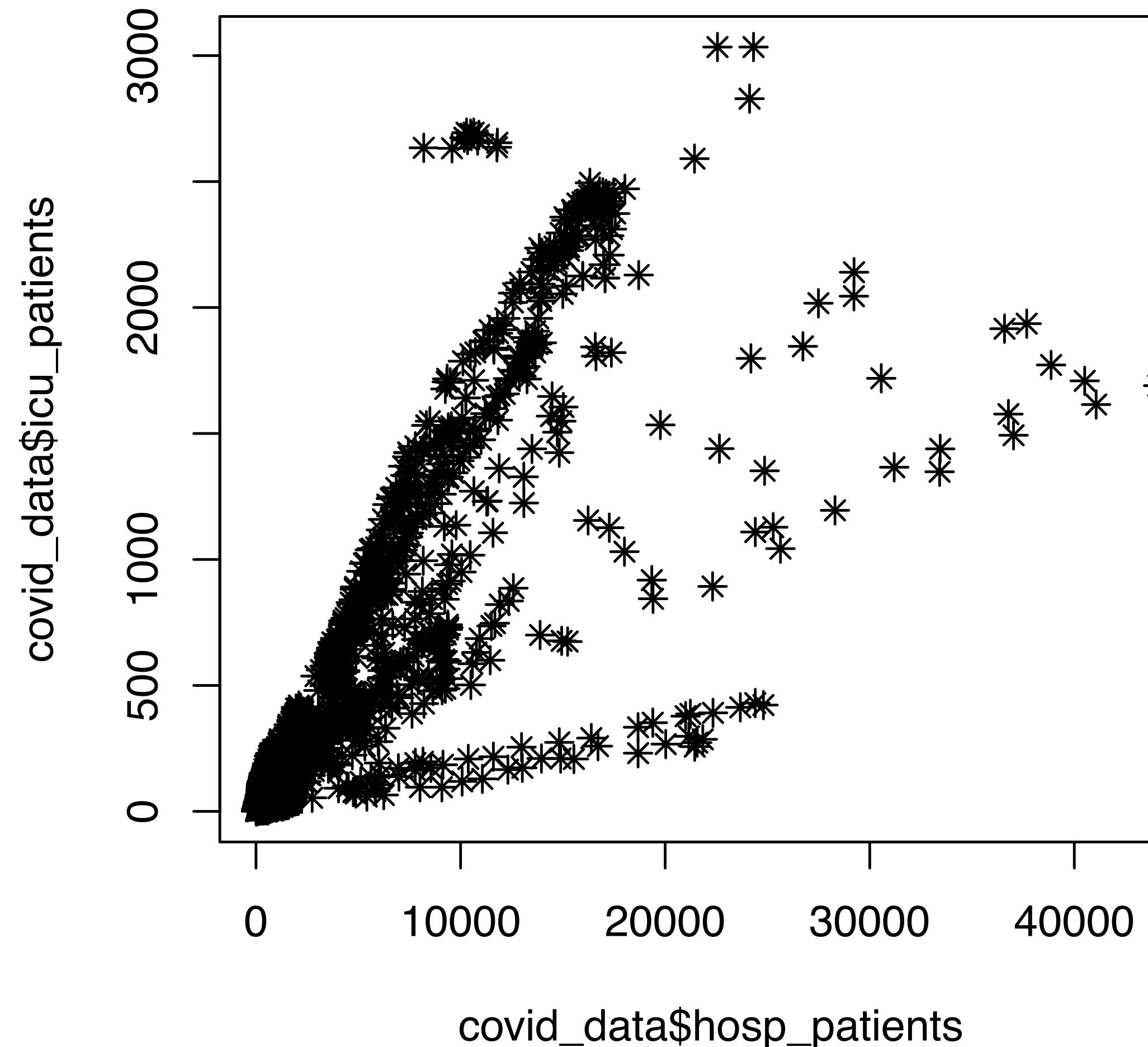
# Anatomy of a scatter plot



data MAPPED TO aesthetic properties OF geometric objects

```
plot(x = covid_data$hosp_patients,  
      y = covid_data$icu_patients,  
      pch = 8,  
      type = "p")
```

# Anatomy of a scatter plot

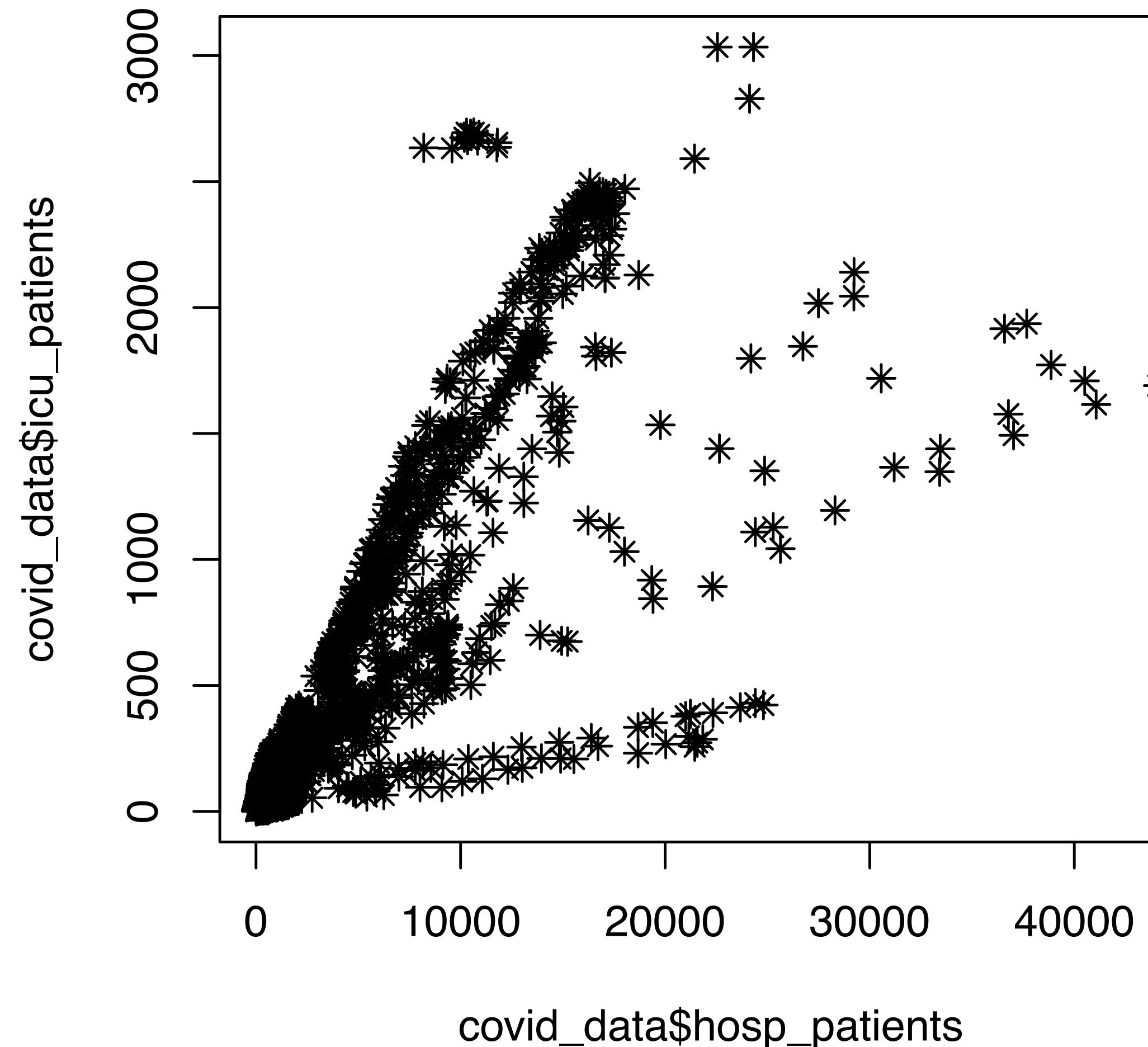


**aesthetic properties**  
OF  
**geometric objects**

**data** MAPPED TO **Data**

```
plot(x = covid_data$hosp_patients,  
      y = covid_data$icu_patients,  
      pch = 8,  
      type = "p")
```

# Anatomy of a scatter plot



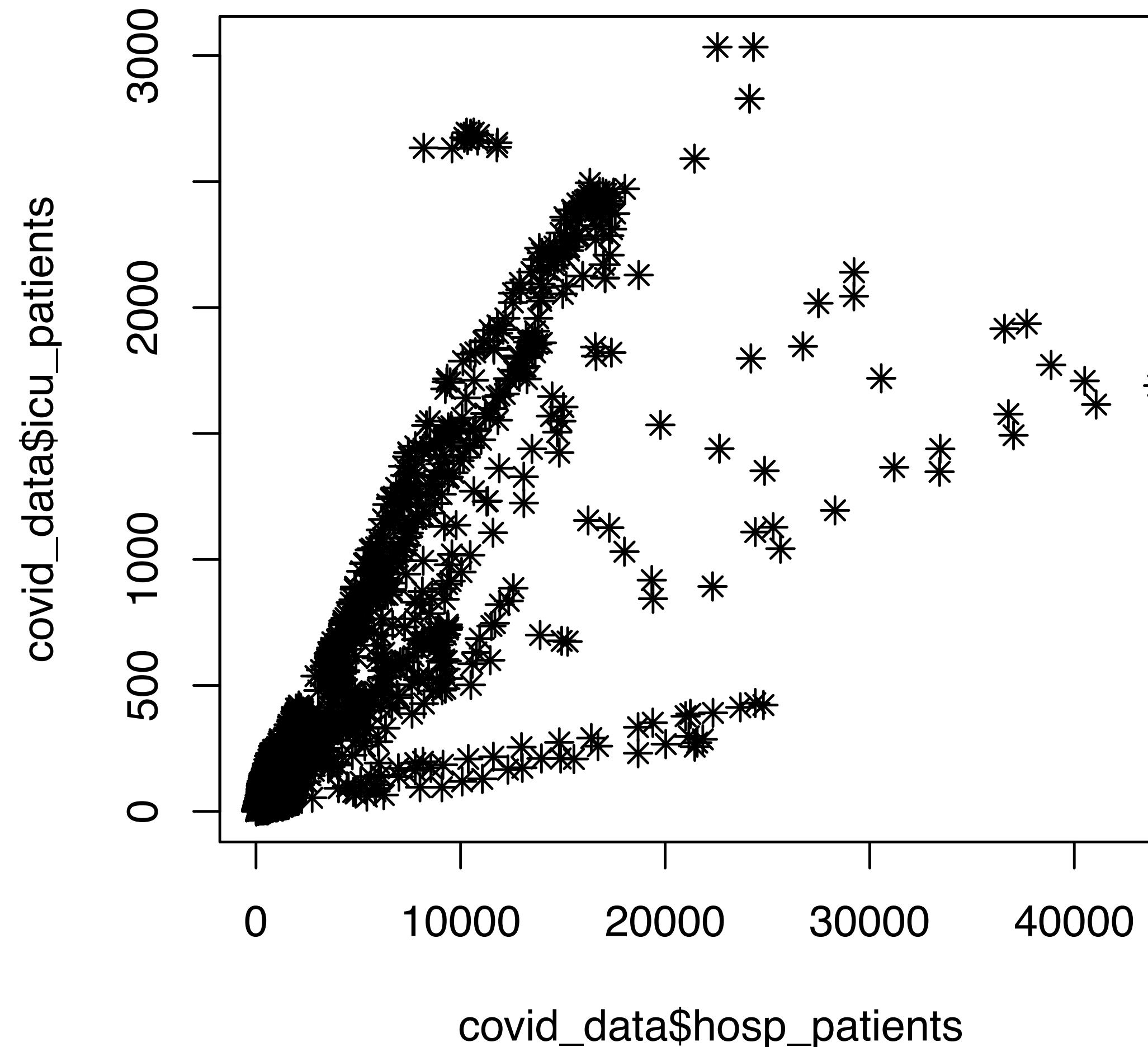
**aesthetic properties** MAPPED TO **geometric objects**

**data** Data

**plot(x = covid\_data\$hosp\_patients,**  
**y = covid\_data\$icu\_patients,**  
**pch = 8,**  
**type = "p")**

**Aesthetic properties**

# Anatomy of a scatter plot



**Variables**

**aesthetic properties**

**of**

**geometric objects**

**MAPPED TO**

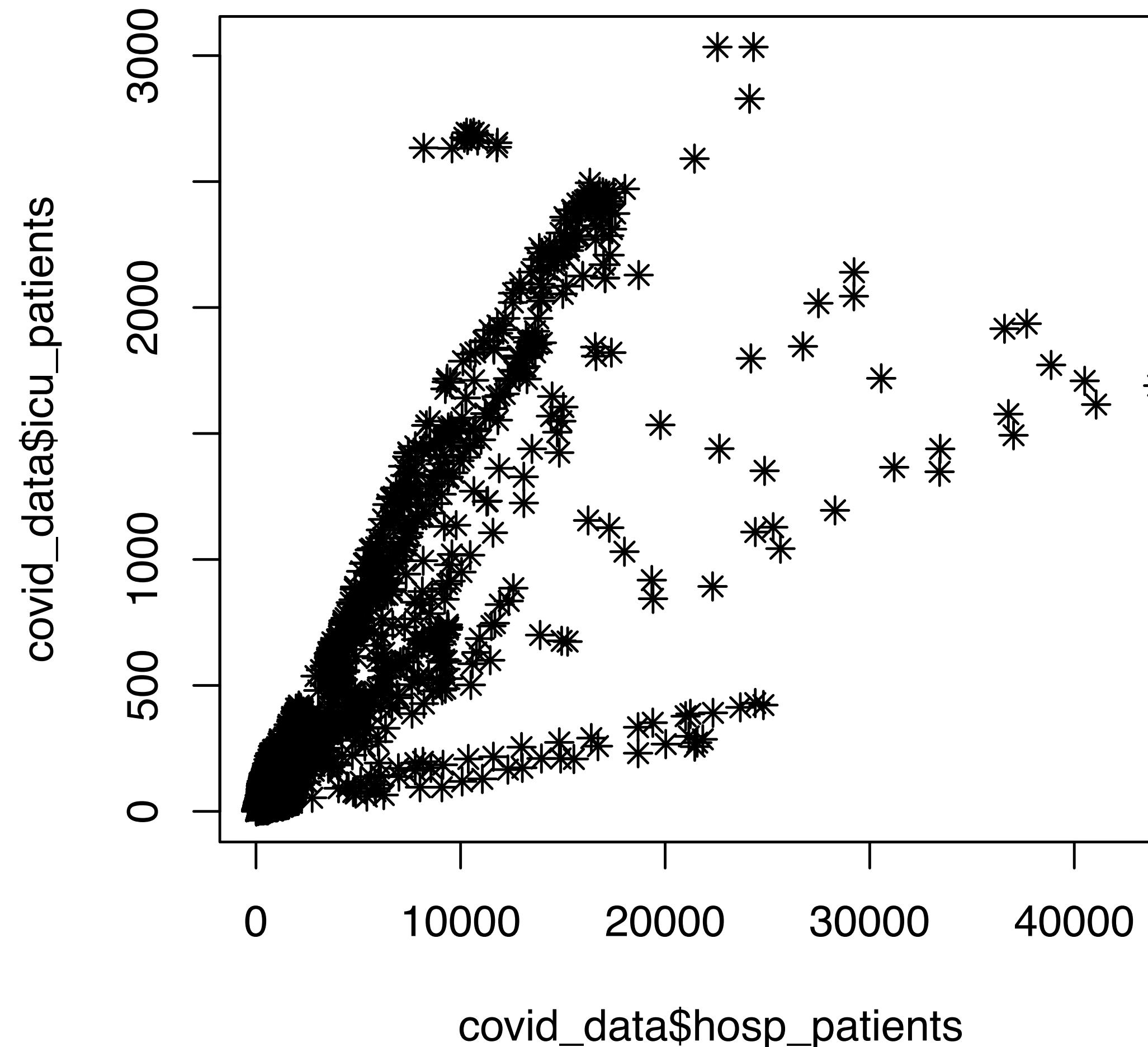
**data**

**Aesthetic properties**

**Data**

**plot(x = covid\_data\$hosp\_patients,**  
**y = covid\_data\$icu\_patients,**  
**pch = 8,**  
**type = “p”)**

# Anatomy of a scatter plot



**aesthetic properties** MAPPED TO **Variables**

**data** Data

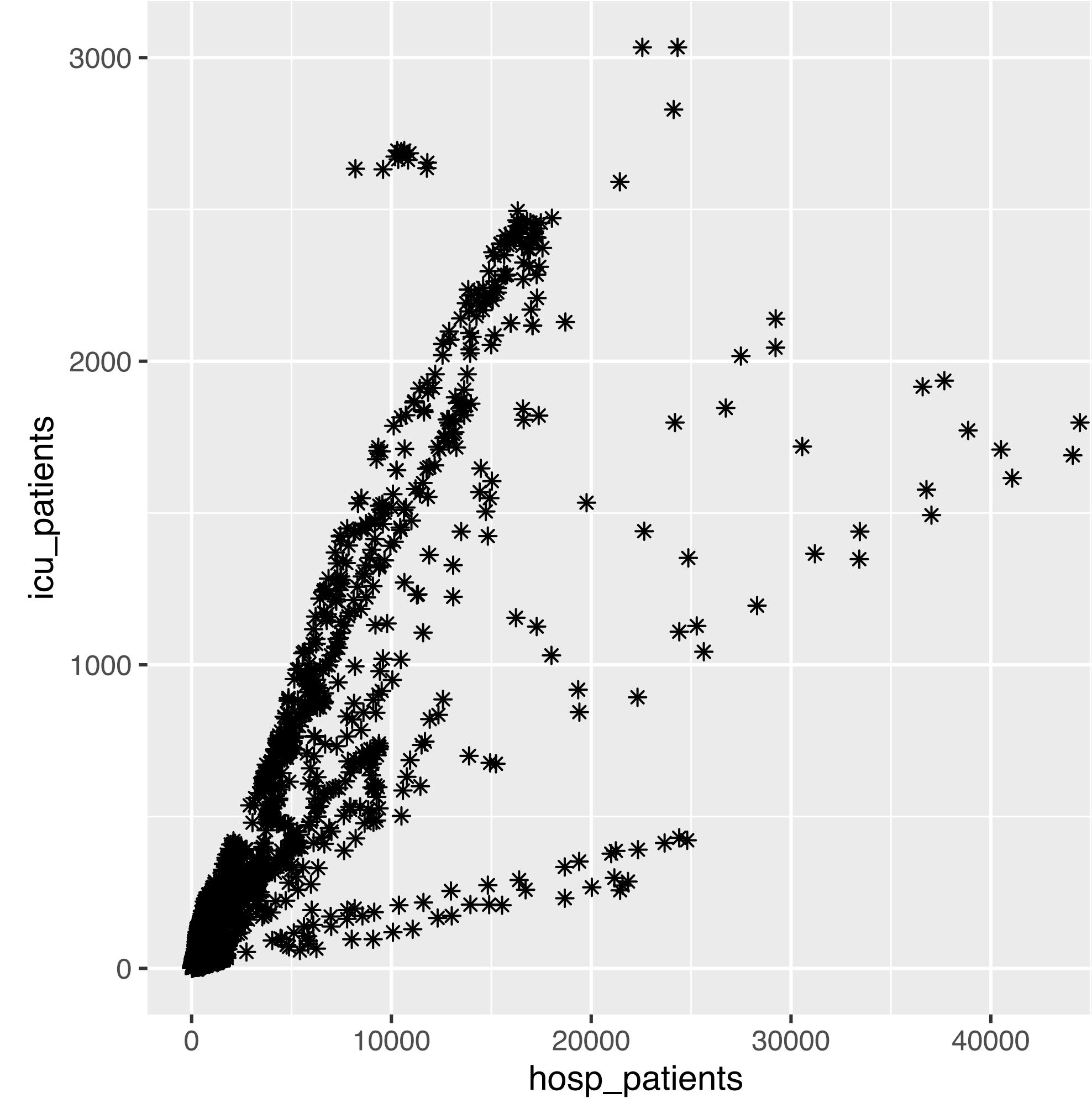
**plot(x = covid\_data\$hosp\_patients,  
 y = covid\_data\$icu\_patients,  
 pch = 8,  
 type = "p")**

**Aesthetic properties**

**Geometric object**

**geometric objects**

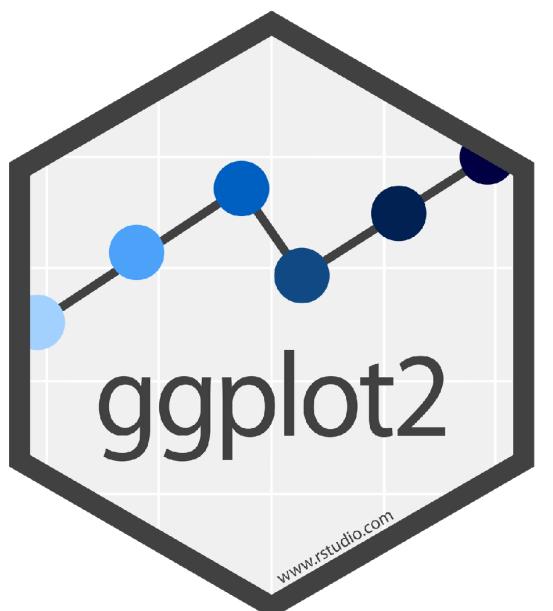
# Anatomy of a simple scatter plot



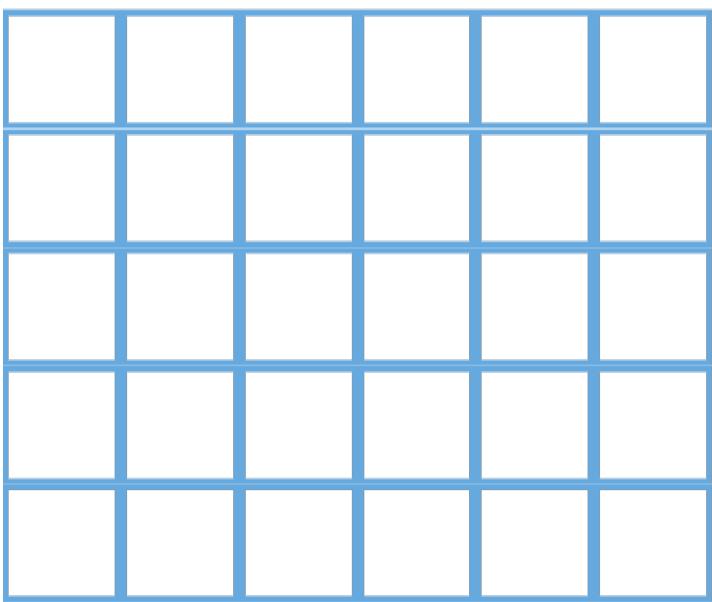
**data** MAPPED TO **aesthetic properties** OF **geometric objects**

```
ggplot(  
  data = covid_data,  
  mapping = aes(  
    x = hosp_patients,  
    y = icu_patients)  
) +  
  geom_point(shape = 8)
```

# A layered grammar of graphics



**data**



MAPPED  
TO

**aesthetic**  
properties

OF

**geometric**  
objects

**x + y** → **geom\_point**

**color** → **geom\_histogram**

**fill** → **geom\_bar**

**alpha** → **geom\_line**

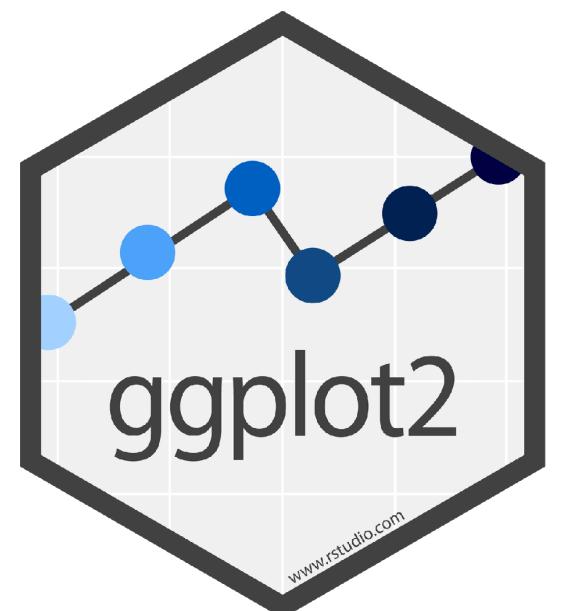
**statistical**  
transformations  
(sometimes)

# Aesthetics

# Geometric shapes

# Scales

# Themes

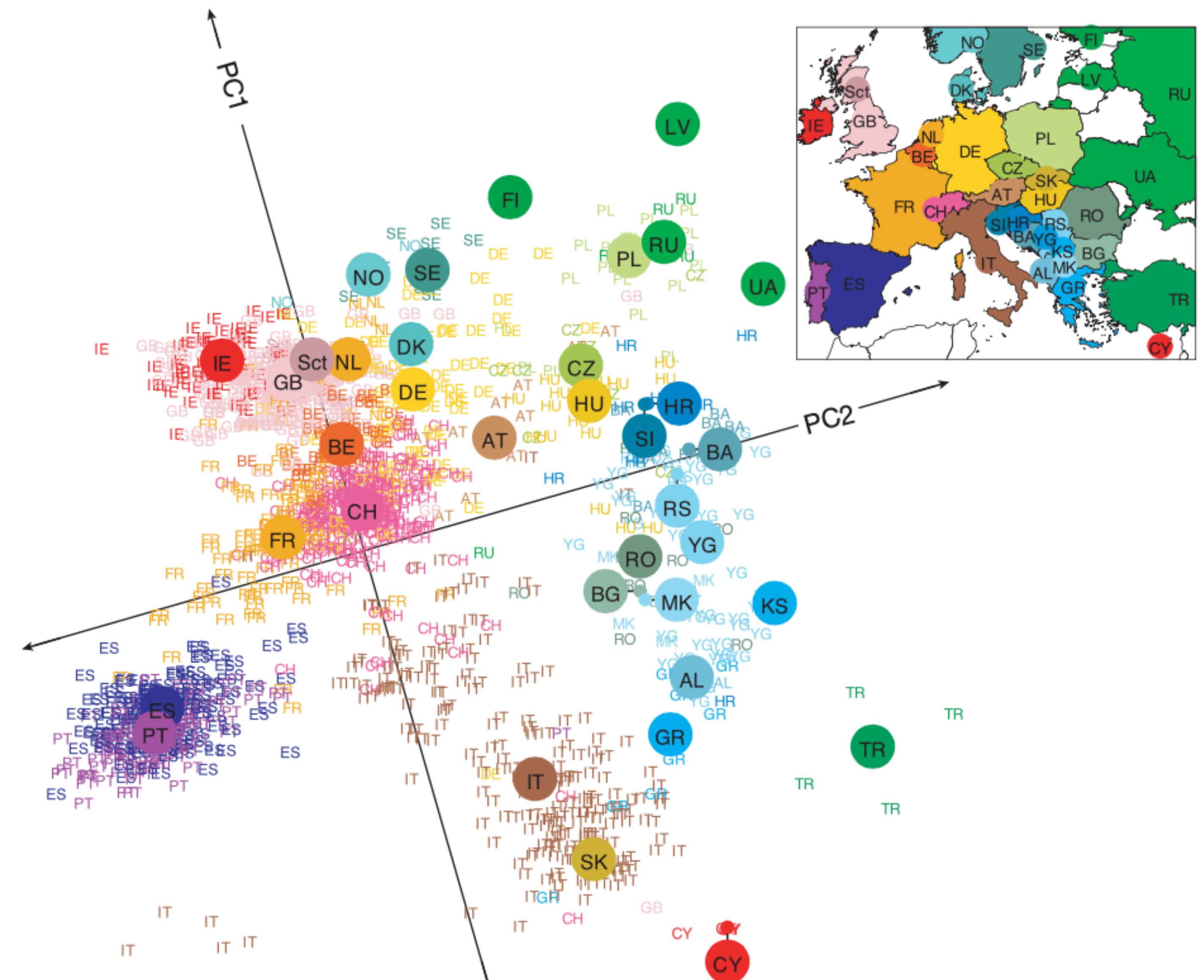
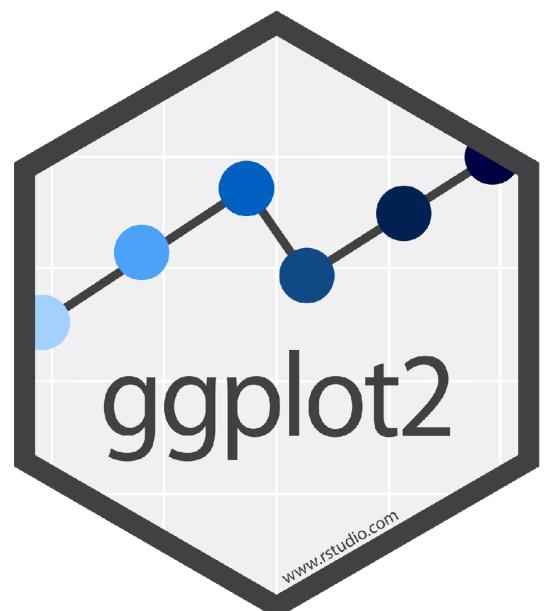


# Aesthetics

## Geometric shapes

## Scales

## Themes



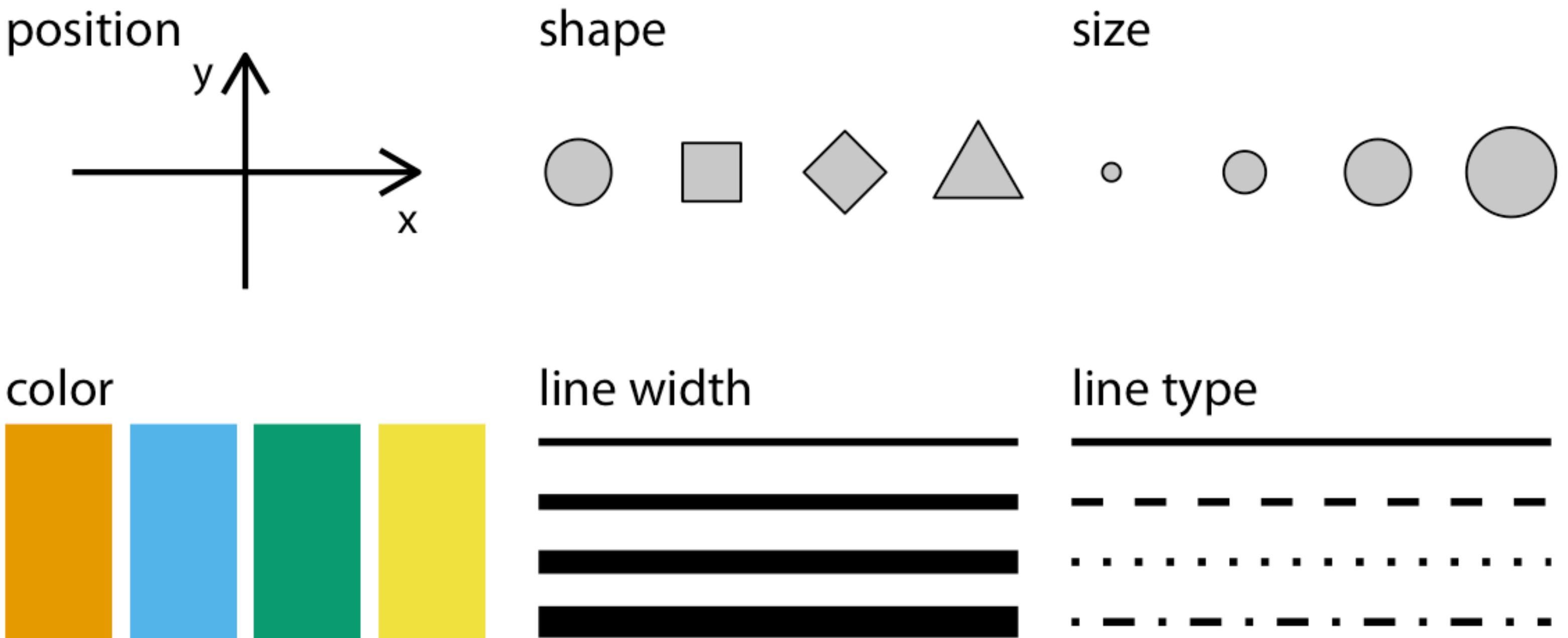
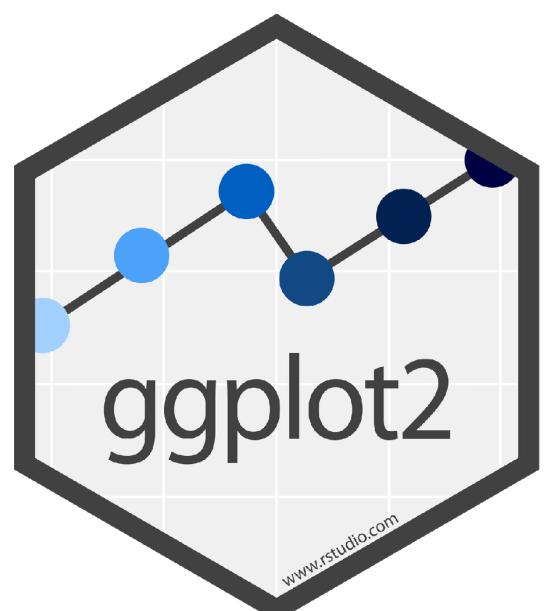
Which aesthetic attributes are used here to visualise attributes of the dataset?

# Aesthetics

## Geometric shapes

## Scales

## Themes



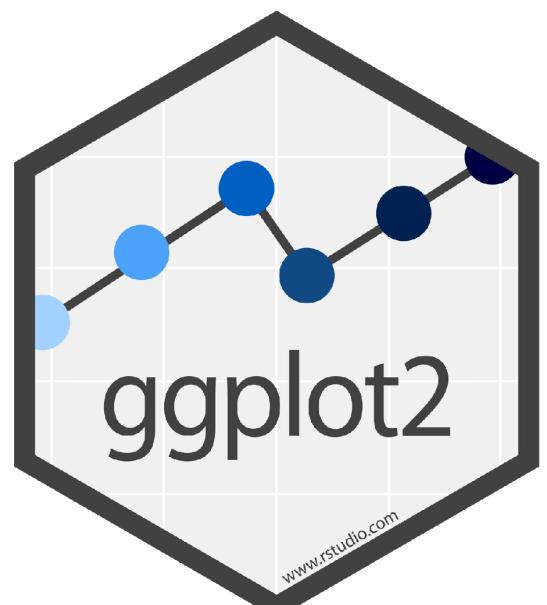
Commonly used aesthetics

# Aesthetics

## Geometric shapes

## Scales

## Themes



## Aes Common aesthetic values.

**color** and **fill** - string ("red", "#RRGGBB")

**linetype** - integer or string (0 = "blank", 1 = "solid", 2 = "dashed", 3 = "dotted", 4 = "dotdash", 5 = "longdash", 6 = "twodash")

**lineend** - string ("round", "butt", or "square")

**linejoin** - string ("round", "mitre", or "bevel")

**size** - integer (line width in mm)

**shape** - integer/shape name or a single character ("a")

0 1 2 3 4 5 6 7 8 9 10 11 12  
□ ○ △ + × ◇ ▽ ■ \* ◆ ⊕ ⊖ ▨  
13 14 15 16 17 18 19 20 21 22 23 24 25  
⊗ □ ○ △ ◇ ○ ○ ● □ ◆ ▲ ▽

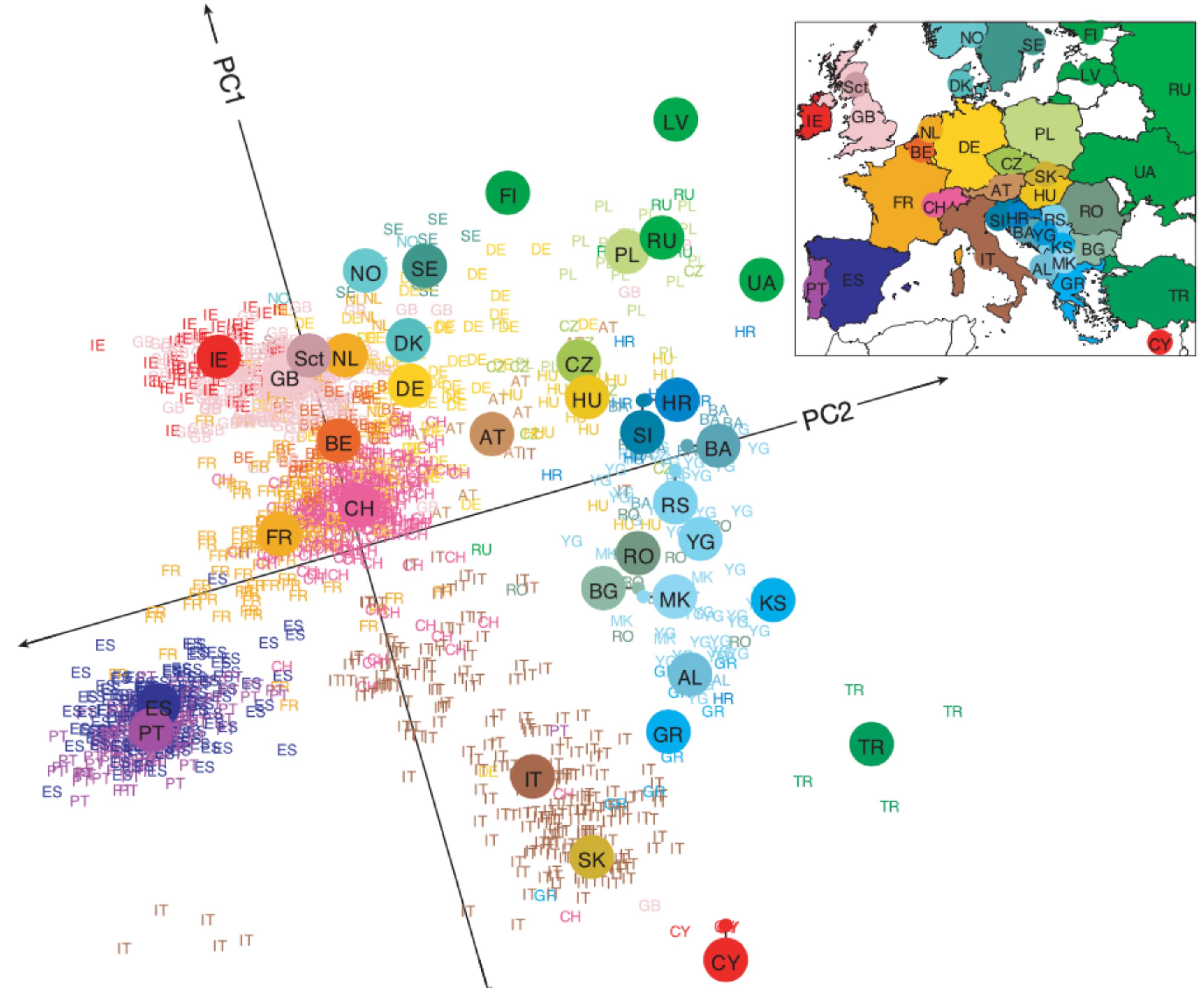
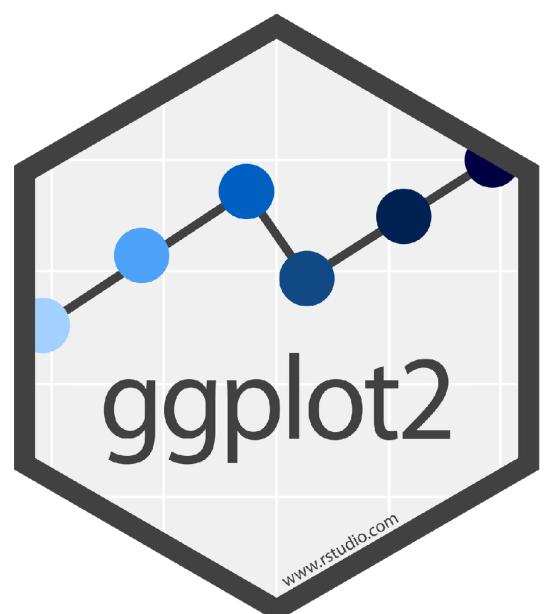
Values for commonly used aesthetics

Aesthetics

## Geometric shapes

Scales

Themes



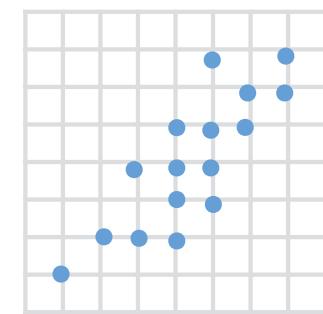
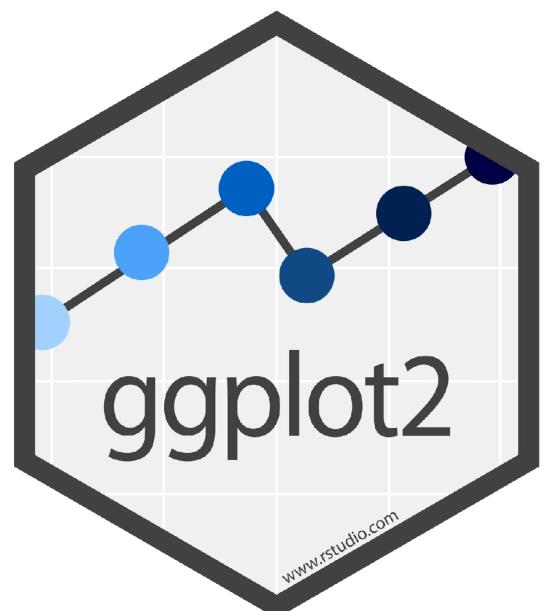
Which geometric shapes are used here to visualise attributes of the dataset?

## Aesthetics

## Geometric shapes

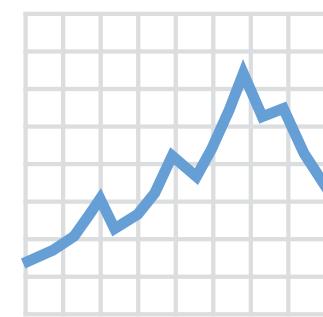
## Scales

## Themes



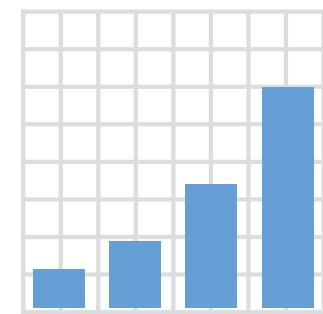
**e + geom\_point()**

x, y, alpha, color, fill, shape, size, stroke



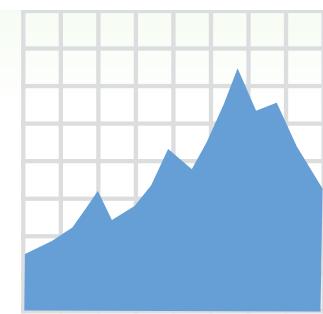
**i + geom\_line()**

x, y, alpha, color, group, linetype, size



**f + geom\_col()**

x, y, alpha, color, fill, group, linetype, size



**c + geom\_area(stat = "bin")**

x, y, alpha, color, fill, linetype, size

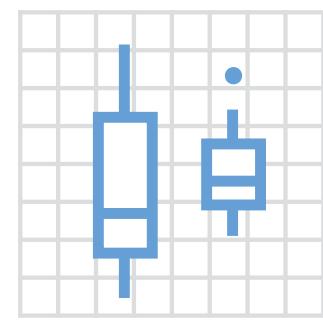
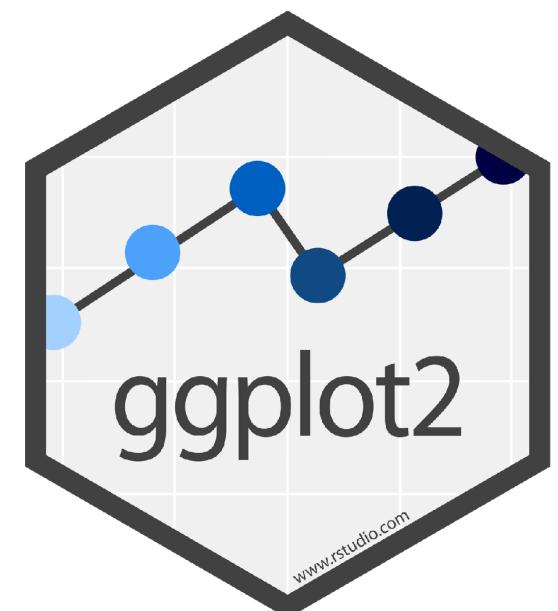
Individual geoms - one distinct object per data point

## Aesthetics

## Geometric shapes

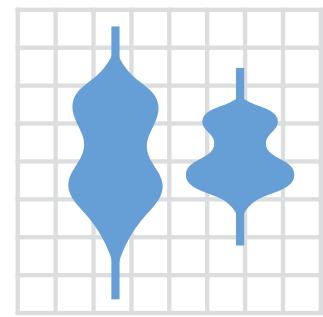
## Scales

## Themes



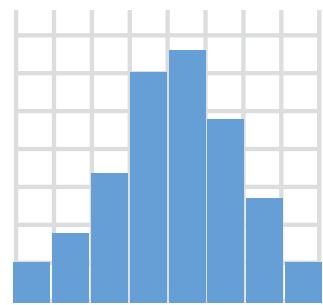
**f + geom\_boxplot()**

x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight



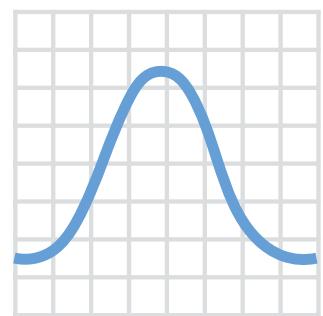
**f + geom\_violin(scale = "area")**

x, y, alpha, color, fill, group, linetype, size, weight



**c + geom\_histogram(binwidth = 5)**

x, y, alpha, color, fill, linetype, size, weight



**c + geom\_density(kernel = "gaussian")**

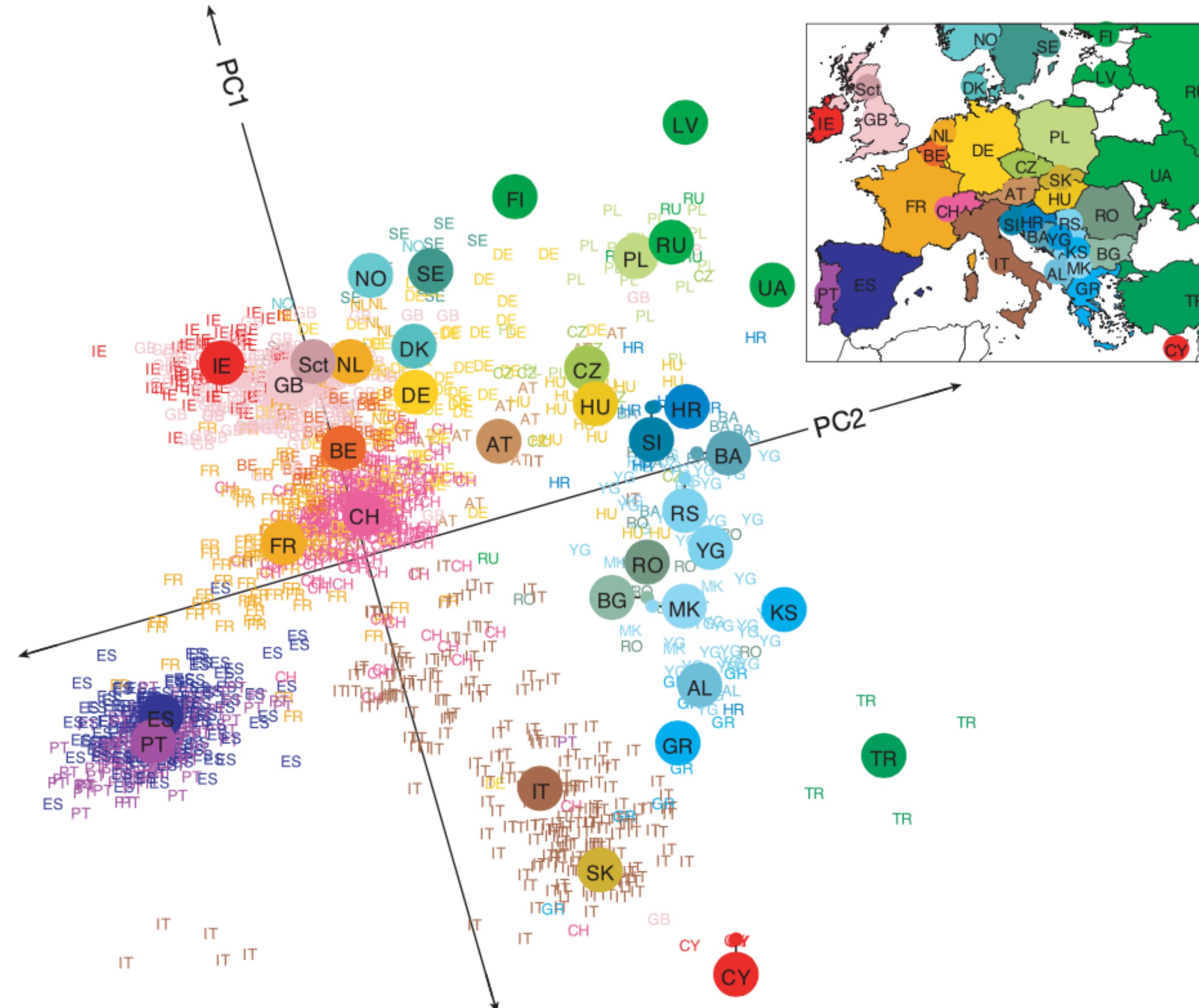
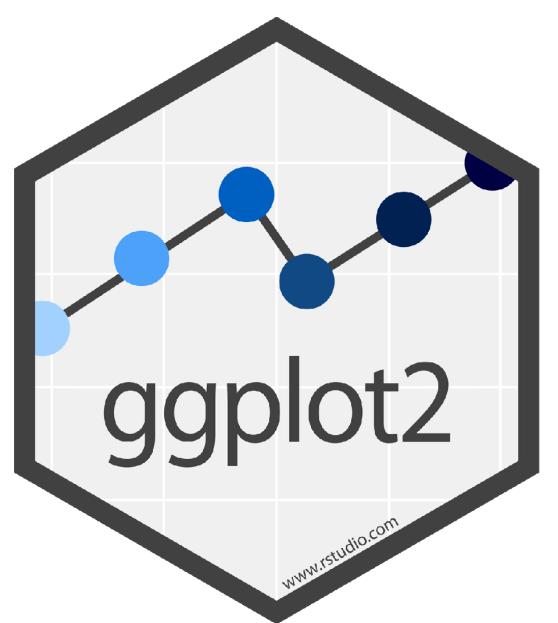
x, y, alpha, color, fill, group, linetype, size, weight

Collective geoms - one distinct object for multiple data points

# Aesthetics

# Geometric shapes

# Scales



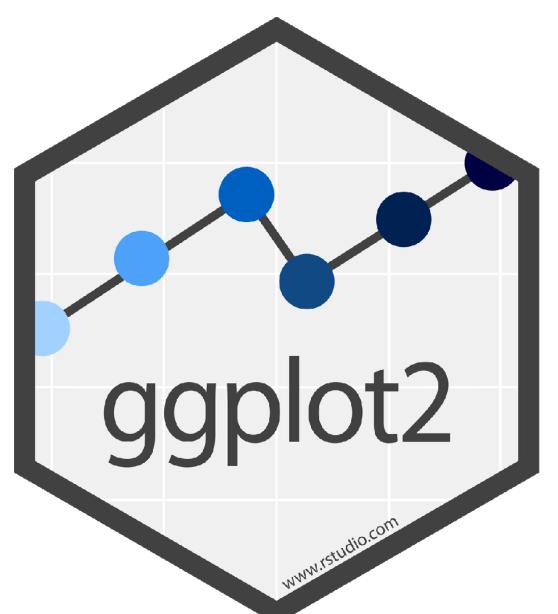
Which scales are used here to visualise attributes of the dataset?

# Aesthetics

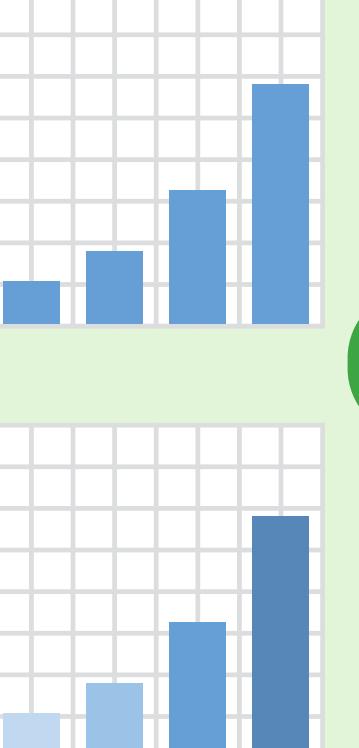
# Geometric shapes

# Scales

# Themes



**Scales** map data values to the visual values of an aesthetic. To change a mapping, add a new scale.



```
n <- d + geom_bar(aes(fill = fl))
```

The code above defines a new layer (`n`) by adding the `geom_bar` layer to the existing data frame (`d`). The `aes(fill = fl)` argument maps the data values to the visual aesthetic of fill color.

Below the code, several green speech bubbles explain the arguments:

- scale\_**: points to the first argument `scale_fill_manual`.
- aesthetic to adjust**: points to the `fill` argument in the `aes` function.
- prepackaged scale to use**: points to the `scale_fill_manual` function itself.
- scale-specific arguments**: points to the three arguments `values`, `limits`, and `name`.
- range of values to include in mapping**: points to the `values` argument.
- title to use in legend/axis**: points to the `limits` argument.
- labels to use in legend/axis**: points to the `name` argument.
- breaks to use in legend/axis**: points to the `limits` argument.

```
n + scale_fill_manual(  
  values = c("skyblue", "royalblue", "blue", "navy"),  
  limits = c("d", "e", "p", "r"), breaks = c("d", "e", "p", "r"),  
  name = "fuel", labels = c("D", "E", "P", "R"))
```

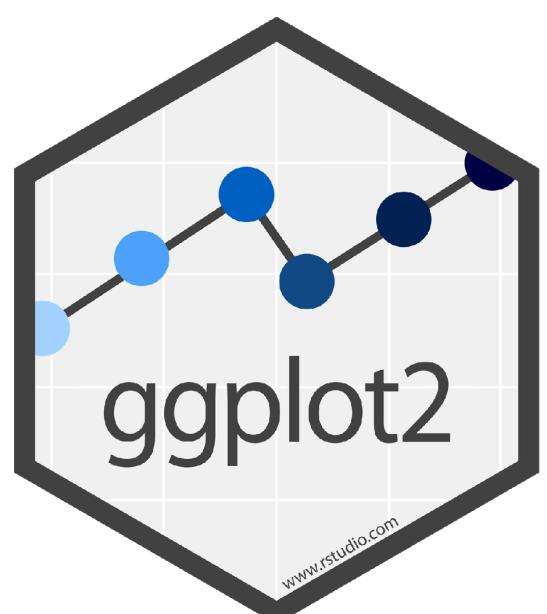
Example of manually setting a scale for aesthetic `fill`

## Aesthetics

## Geometric shapes

## Scales

## Themes



**scale\_\*\_continuous()** - Map cont' values to visual ones.  
**scale\_\*\_discrete()** - Map discrete values to visual ones.  
**scale\_\*\_binned()** - Map continuous values to discrete bins.  
**scale\_\*\_identity()** - Use data values as visual ones.  
**scale\_\*\_manual(values = c())** - Map discrete values to manually chosen visual ones.

**scale\_x\_log10()** - Plot x on log10 scale.  
**scale\_x\_reverse()** - Reverse the direction of the x axis.  
**scale\_x\_sqrt()** - Plot x on square root scale.

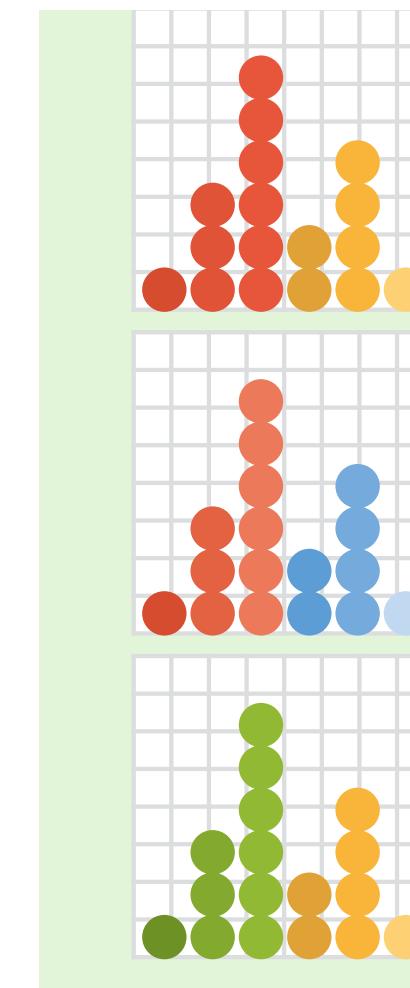
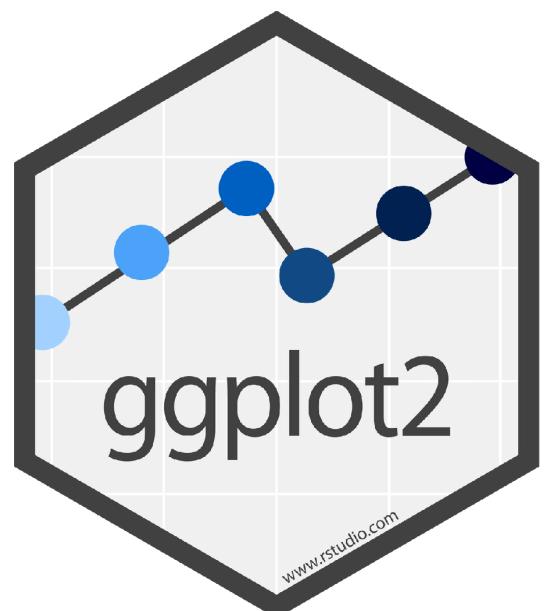
Examples of different scales

# Aesthetics

## Geometric shapes

## Scales

## Themes

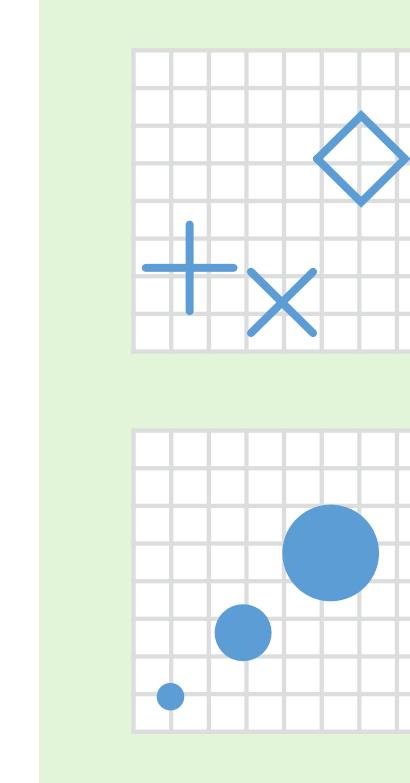


**`o + scale_fill_gradient(low="red", high="yellow")`**

**`o + scale_fill_gradient2(low = "red", high = "blue", mid = "white", midpoint = 25)`**

**`o + scale_fill_gradientn(colors = topo.colors(6))`**

Also: `rainbow()`, `heat.colors()`, `terrain.colors()`,  
`cm.colors()`, `RColorBrewer::brewer.pal()`



**`p + scale_shape() + scale_size()`**

**`p + scale_shape_manual(values = c(3:7))`**

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25

□○△+×◊▽⊗✳⊕⊛⊛田⊗□○△◊○○●□◆△▽

**`p + scale_radius(range = c(1,6))`**

**`p + scale_size_area(max_size = 6)`**

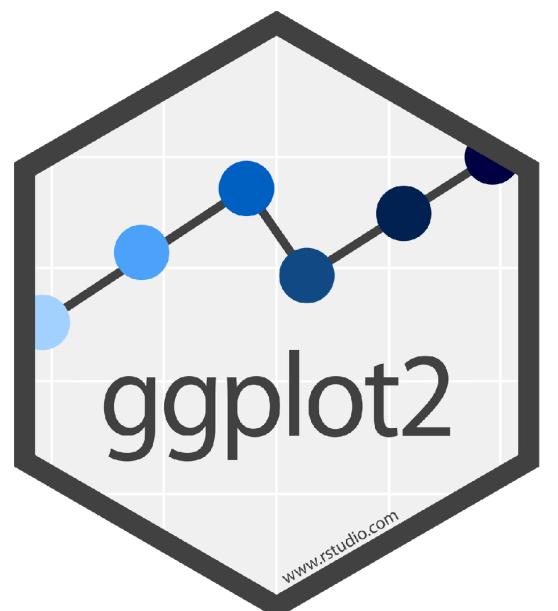
Examples of different scales

## Aesthetics

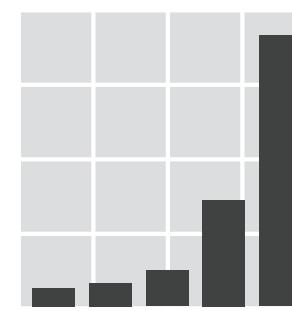
## Geometric shapes

## Scales

## Themes



**r + theme()** Customize aspects of the theme such as axis, legend, panel, and facet properties.



**r + theme\_gray()**  
Grey background  
(default theme).

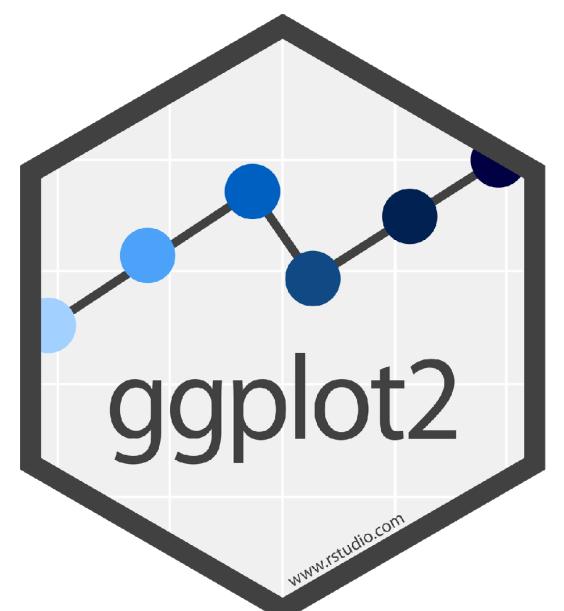
Themes control the non-data elements of a plot

## Aesthetics

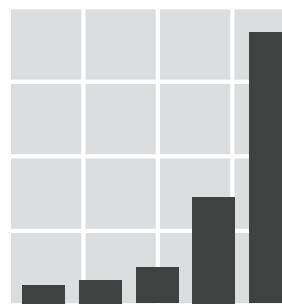
## Geometric shapes

## Scales

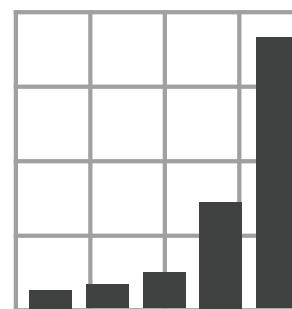
## Themes



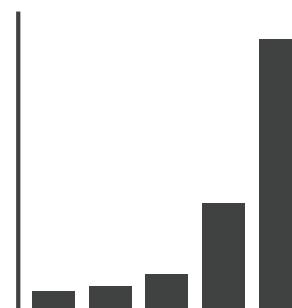
**r + theme()** Customize aspects of the theme such as axis, legend, panel, and facet properties.



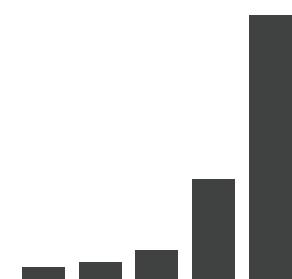
**r + theme\_gray()**  
Grey background  
(default theme).



**r + theme\_bw()**  
White background  
with grid lines.



**r + theme\_classic()**



**r + theme\_void()**  
Empty theme.

Examples of available pre-defined themes

**[https://ucph.padlet.org/martinsikora4/data\\_analysis\\_2025\\_dataviz](https://ucph.padlet.org/martinsikora4/data_analysis_2025_dataviz)**