

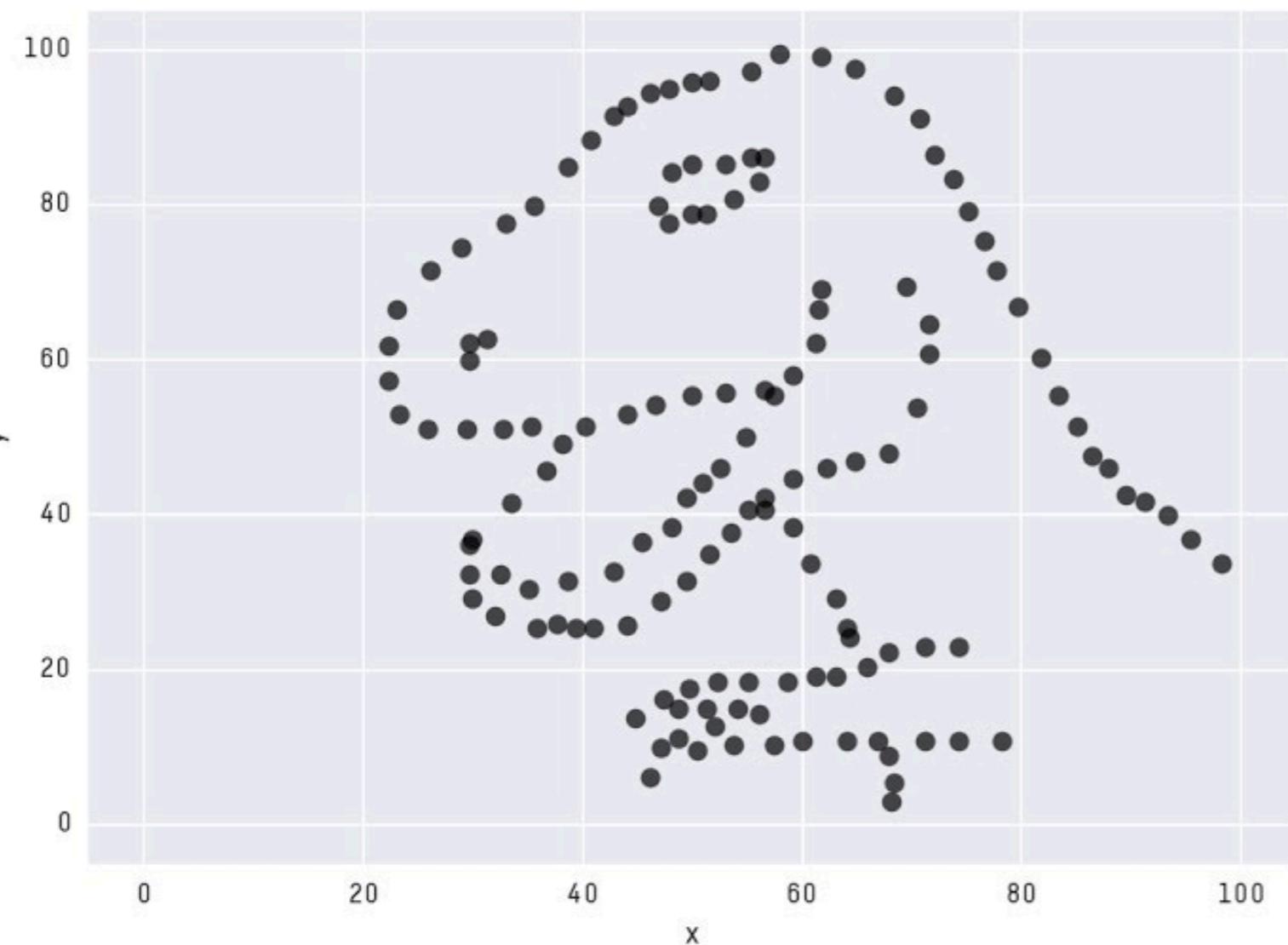
Data Vis with ggplot2

Why and How

<http://bit.ly/2xEUGIu>

Przemysław Biecek
<http://biecek.pl>

Package datasauRus



X Mean : 54.2632025
Y Mean : 47.8315781
X SD : 16.7650109
Y SD : 26.9353144
Corr. : -0.0645195

Grammar of Graphics

ggplot2

Three ecosystems for static statistical graphics

```
library(PBImisc)
```

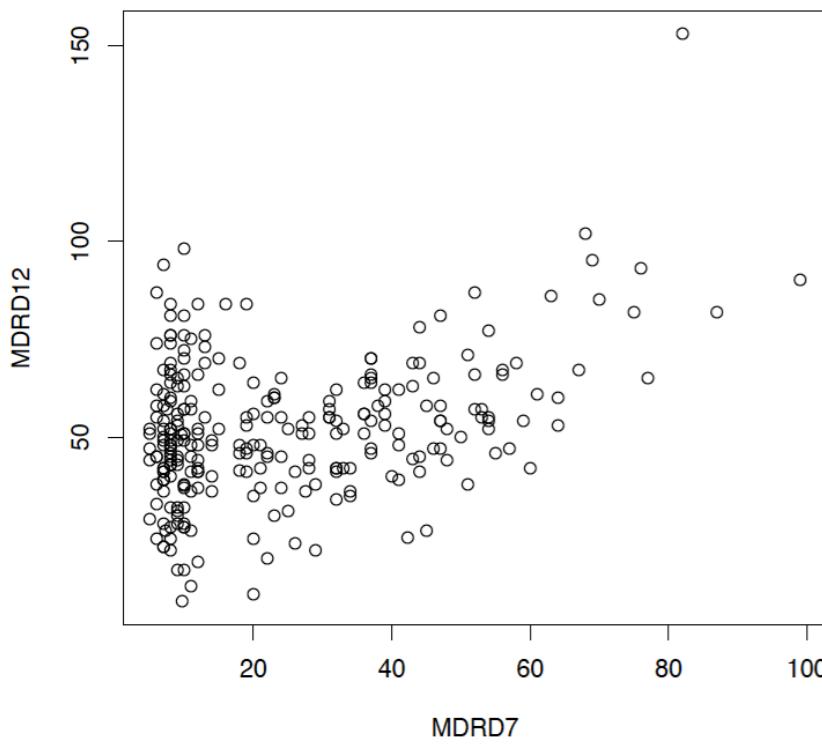
```
# library graphics  
plot(MDRD12~MDRD7, kidney)
```

```
# library lattice  
xyplot(MDRD12~MDRD7, kidney)
```

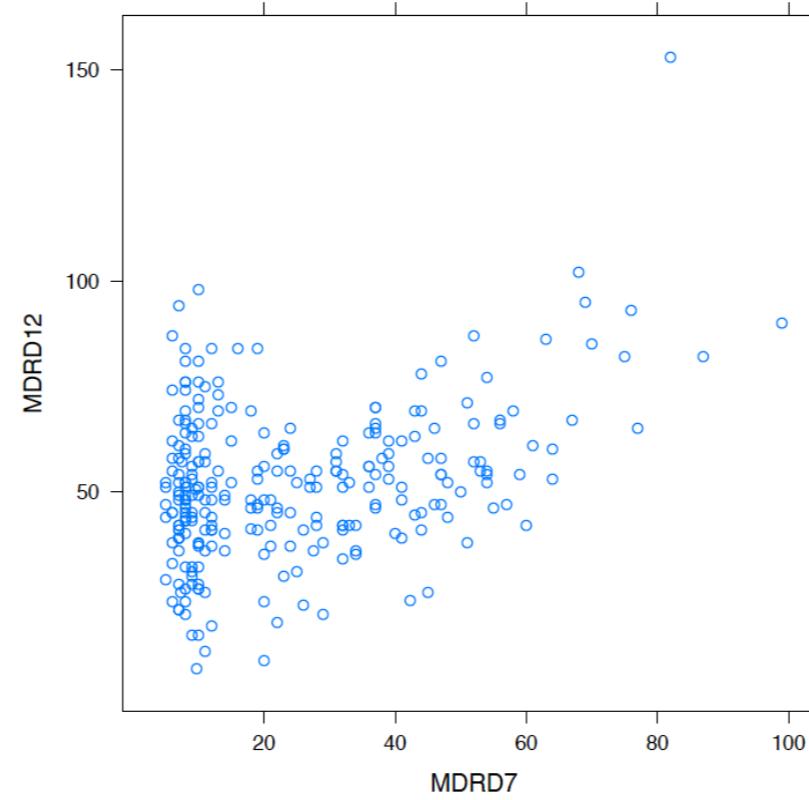
```
# library ggplot2  
qplot(MDRD12, MDRD7, data=kidney)
```

Find differences between these plots

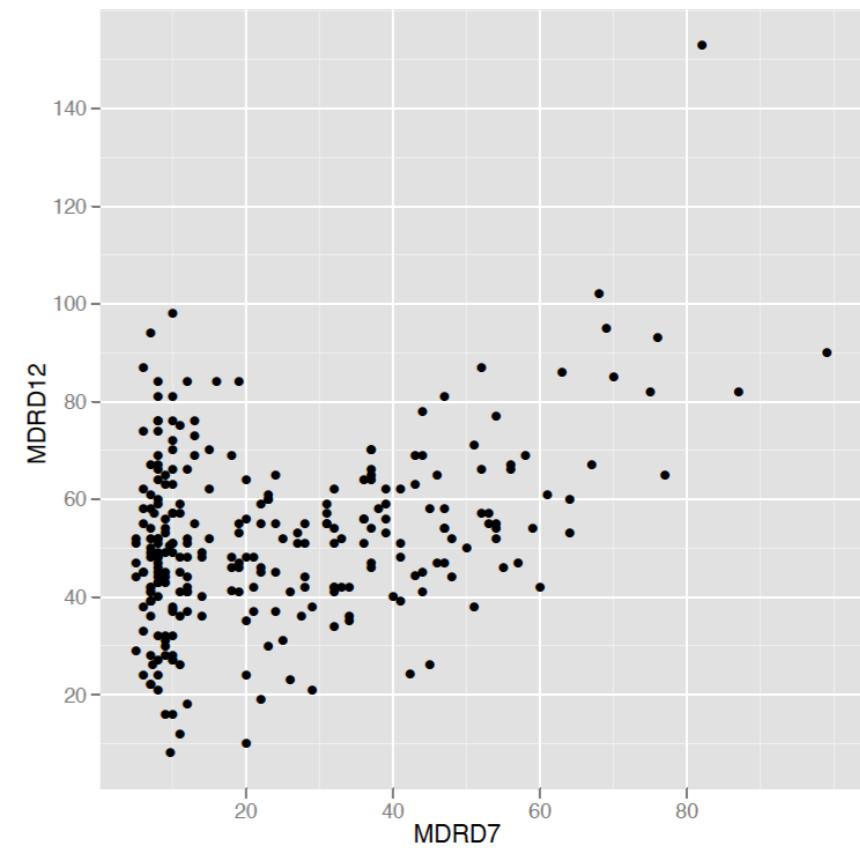
graphics



lattice



ggplot2

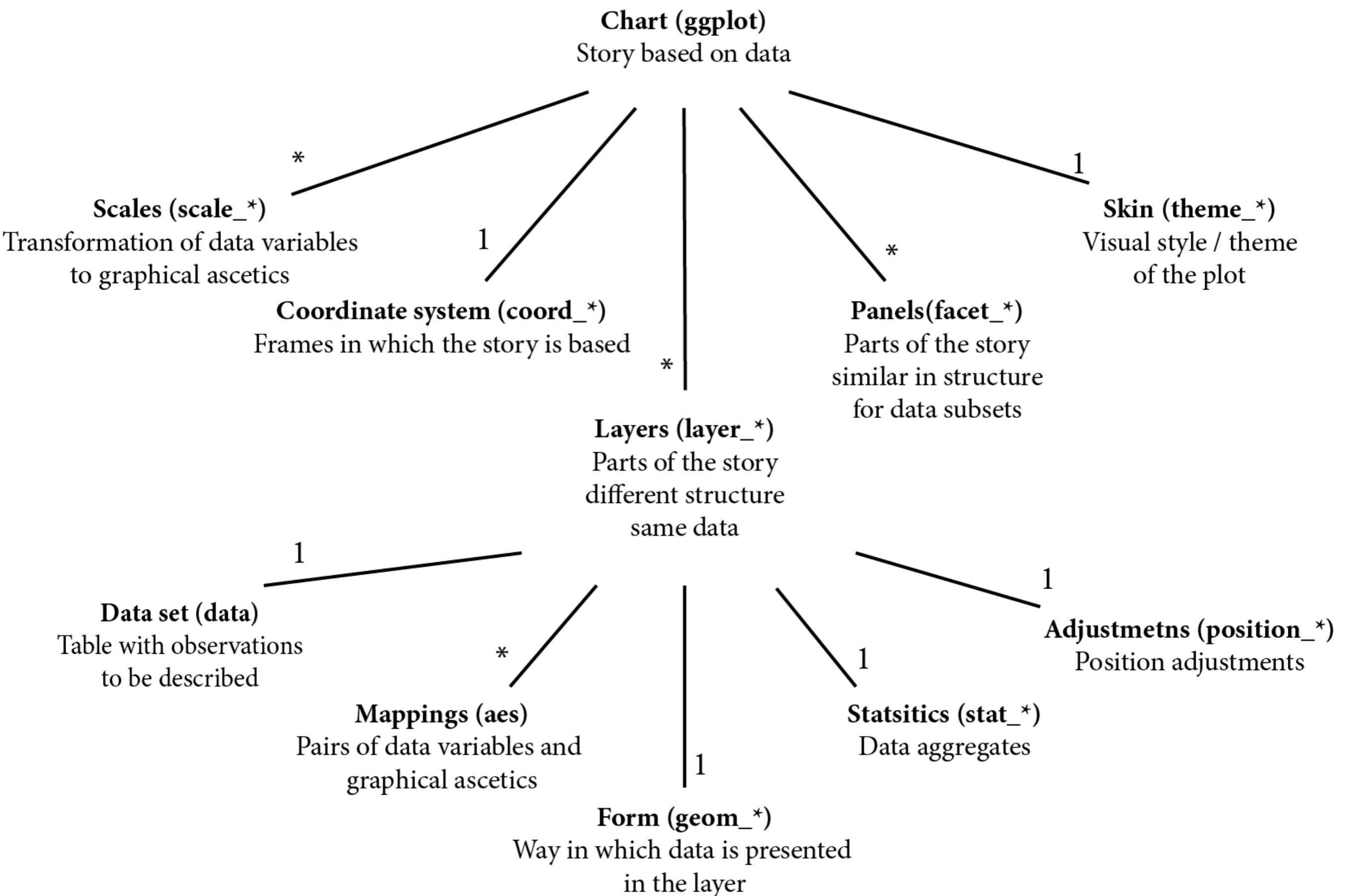


Why ggplot2?

- + Elegant
- + Highly customisable
- + Uniform
- + Natural
- + Expressive
- + Popular
- Steep learning curve
- Slow
- Evolving pretty fast
(too fast?)

Why ggplot2?





[hands-on live R]

```
# Best 250 series http://www.imdb.com/chart/toptv/
## (1) read the new data with archivist
library(archivist)
series2017 <- aread("mi2-warsaw/RLadies/arepo/45aa16dc4dbf0d87e3e40eb9dc9d18ae")

## (2) or read the old data with Pogromcy Danych
library(PogromcyDanych)
serialIMDB

## (3) or scrap the data from IMDB database
library(rvest)
library(dplyr)
# read links and titles
page <- read_html("http://www.imdb.com/chart/toptv/")
series <- html_nodes(page, ".titleColumn a")
titles <- html_text(series)
links <- html_attr(series, "href")
codes <- sapply(strsplit(links, split = "/"), `[,` , 3)

allseries <- lapply(seq_along(codes), function(i) {
  tab <- read_html(paste0("http://www.imdb.com/title/", codes[i], "/epdate?ref_=ttep_q1_4")) %>%
    html_node("table") %>%
    html_table()
  data.frame(Serie = titles[i], tab[,1:4], Season = gsub(tab[,1], pattern="\.\.*", replacement=""))
})
```

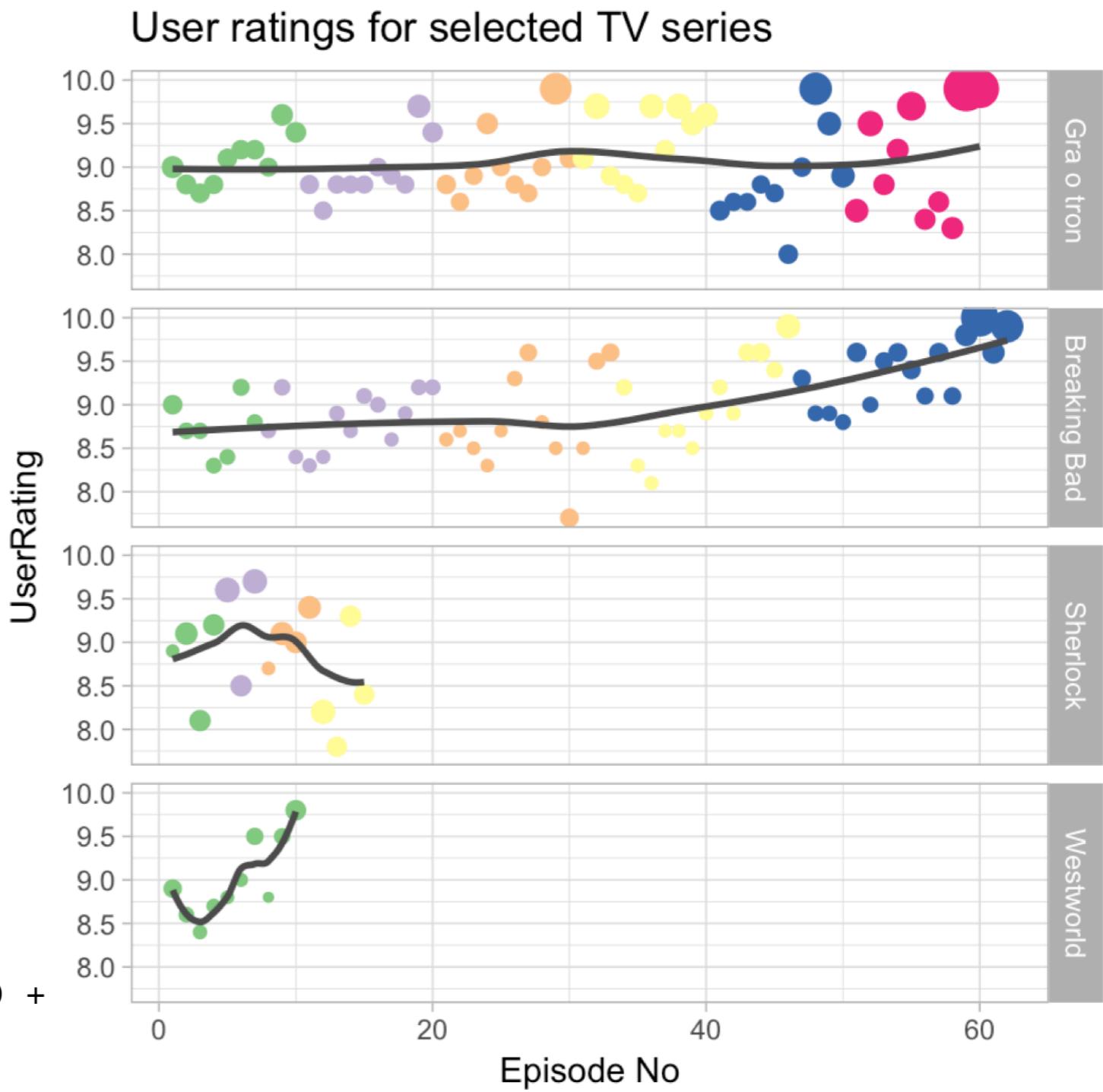
```

selected <- c("Game of Thrones", "Breaking Bad",
           "Sherlock", "Westworld")

dat <- series2017 %>%
  filter(Serie %in% selected)

ggplot(dat, aes(id, UserRating)) +
  geom_point(aes(color=Season, size=UserVotes)) +
  geom_smooth(se=FALSE, color="grey30") +
  facet_grid(Serie~.) +
  theme_light() + theme(legend.position="none") +
  scale_color_brewer(palette = 1, type = "qual") +
  ggtitle("User ratings for selected TV series") +
  xlab("Episode No")

```



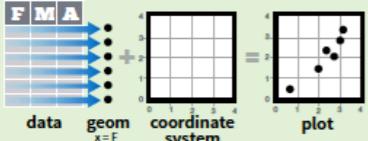
Data Visualization with ggplot2

Cheat Sheet

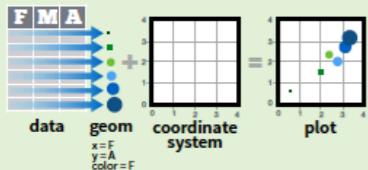


Basics

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()**

aesthetic mappings **data** **geom**
`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
`ggplot(mpg, aes(hwy, cty)) + geom_point(aes(color = cyl)) + geom_smooth(method = "lm") + coord_cartesian() + scale_color_gradient() + theme_bw()`
add layers, elements with +
layer = geom + default stat + layer specific 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)`

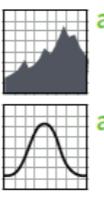
Saves last plot as 5' x 5' file named "plot.png" in working directory. Matches file type to file extension.

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 <- ggplot(mpg, aes(hwy))`



a + geom_area(stat = "bin")

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

b + geom_area(aes(y = ..density..), stat = "bin")

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

a + geom_density(kernel = "gaussian")

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

b + geom_density(aes(y = ..county..))

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

a + geom_dotplot()

x, y, alpha, color, fill

a + geom_freqpoly()

x, y, alpha, color, linetype, size

b + geom_freqpoly(aes(y = ..density..))

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

a + geom_histogram(binwidth = 5)

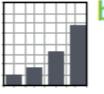
x, y, alpha, color, fill, linetype, size, weight

b + geom_histogram(aes(y = ..density..))

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

Discrete

`b <- ggplot(mpg, aes(fl))`



b + geom_bar()

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

Graphical Primitives

`c <- ggplot(map, aes(long, lat))`



c + geom_polygon(aes(group = group))

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

`d <- ggplot(economics, aes(date, unemploy))`



d + geom_path(lineend = "butt",

linejoin = "round", linemitre = 1)

x, y, alpha, color, linetype, size

d + geom_ribbon(aes(ymin = unemploy - 900,

ymax = unemploy + 900))

x, ymax, ymin, alpha, color, fill, linetype, size

`e <- ggplot(seals, aes(x = long, y = lat))`



e + geom_segment(aes(

xend = long + delta_long,

yend = lat + delta_lat))

x, xend, y, yend, alpha, color, linetype, size

e + geom_rect(aes(xmin = long, ymin = lat,

xmax = long + delta_long,

ymax = lat + delta_lat))

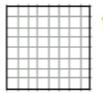
xmax, xmin, ymax, ymin, alpha, color, fill,

linetype, size

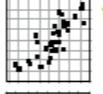
Two Variables

Continuous X, Continuous Y

`f <- ggplot(mpg, aes(cty, hwy))`

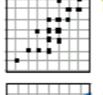


f + geom_blank()



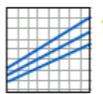
f + geom_jitter()

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



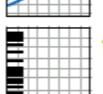
f + geom_point()

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



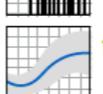
f + geom_quantile()

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



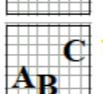
f + geom_rug(sides = "bl")

alpha, color, linetype, size



f + geom_smooth(model = lm)

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



f + geom_text(aes(label = cty))

x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

Three Variables

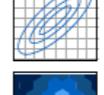
Continuous Bivariate Distribution

`i <- ggplot(movies, aes(year, rating))`



i + geom_bin2d(binwidth = c(5, 0.5))

xmax, xmin, ymax, ymin, alpha, color, fill, linetype, size, weight



i + geom_density2d()

x, y, alpha, colour, linetype, size

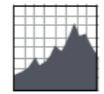


i + geom_hex()

x, y, alpha, colour, fill size

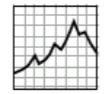
Continuous Function

`j <- ggplot(economics, aes(date, unemploy))`



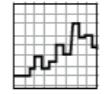
j + geom_area()

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



j + geom_line()

x, y, alpha, color, linetype, size



j + geom_step(direction = "hv")

x, y, alpha, color, linetype, size

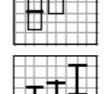
Visualizing error

`df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)`



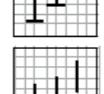
k + geom_crossbar(fatten = 2)

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



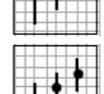
k + geom_errorbar()

x, ymax, ymin, alpha, color, linetype, size, width (also **geom_errorbarh()**)



k + geom_linerange()

x, ymin, ymax, alpha, color, linetype, size



k + geom_pointrange()

x, y, ymin, ymax, alpha, color, fill, linetype, shape, size

Maps

`data <- data.frame(murder = USArrests$Murder, state = tolower(rownames(USArrests)))`

`map <- map_data("state")`



l + geom_map(aes(map_id = state), map = map) +

expand_limits(x = map\$long, y = map\$lat)

map_id, alpha, color, fill, linetype, size

Three Variables

`seals$z <- with(seals, sqrt(delta_long^2 + delta_lat^2))`

`m <- ggplot(seals, aes(long, lat))`



m + geom_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)

x, y, alpha, fill



m + geom_contour(aes(z = z))

x, y, z, alpha, colour, linetype, size, weight

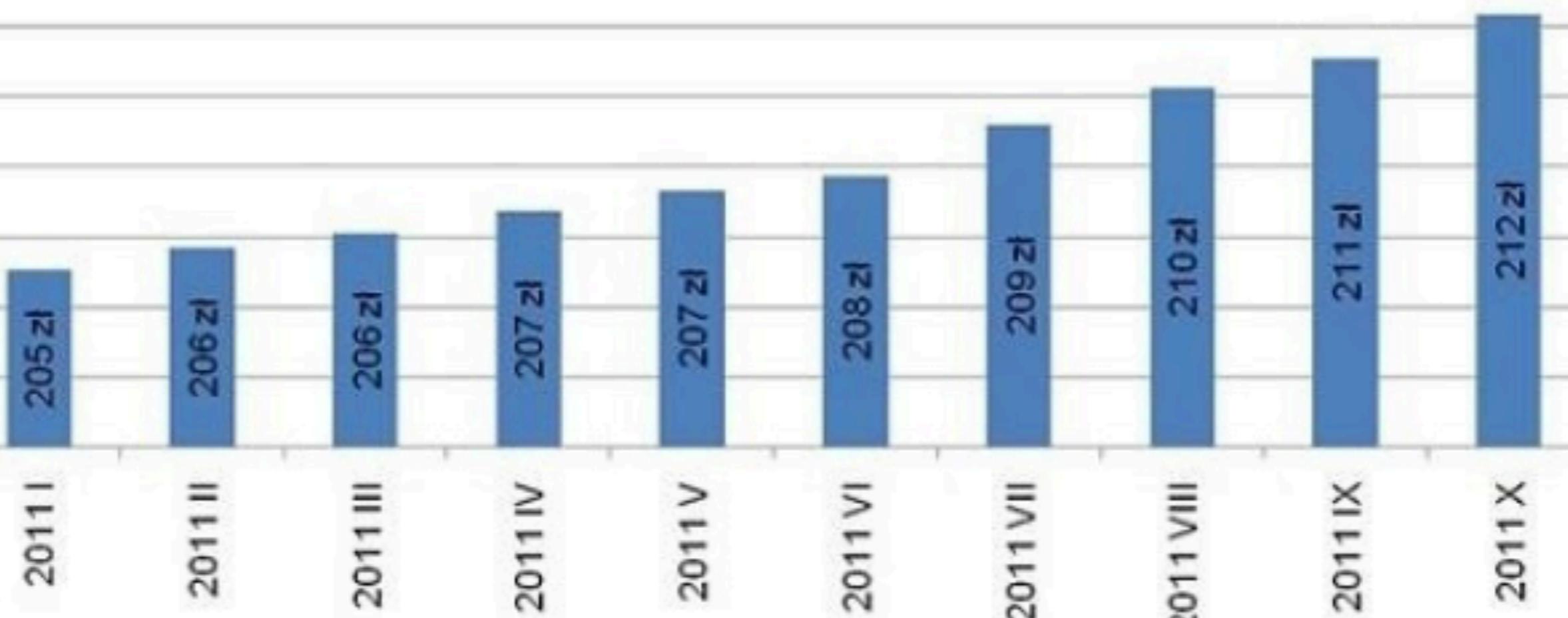


m + geom_tile(aes(fill = z))

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

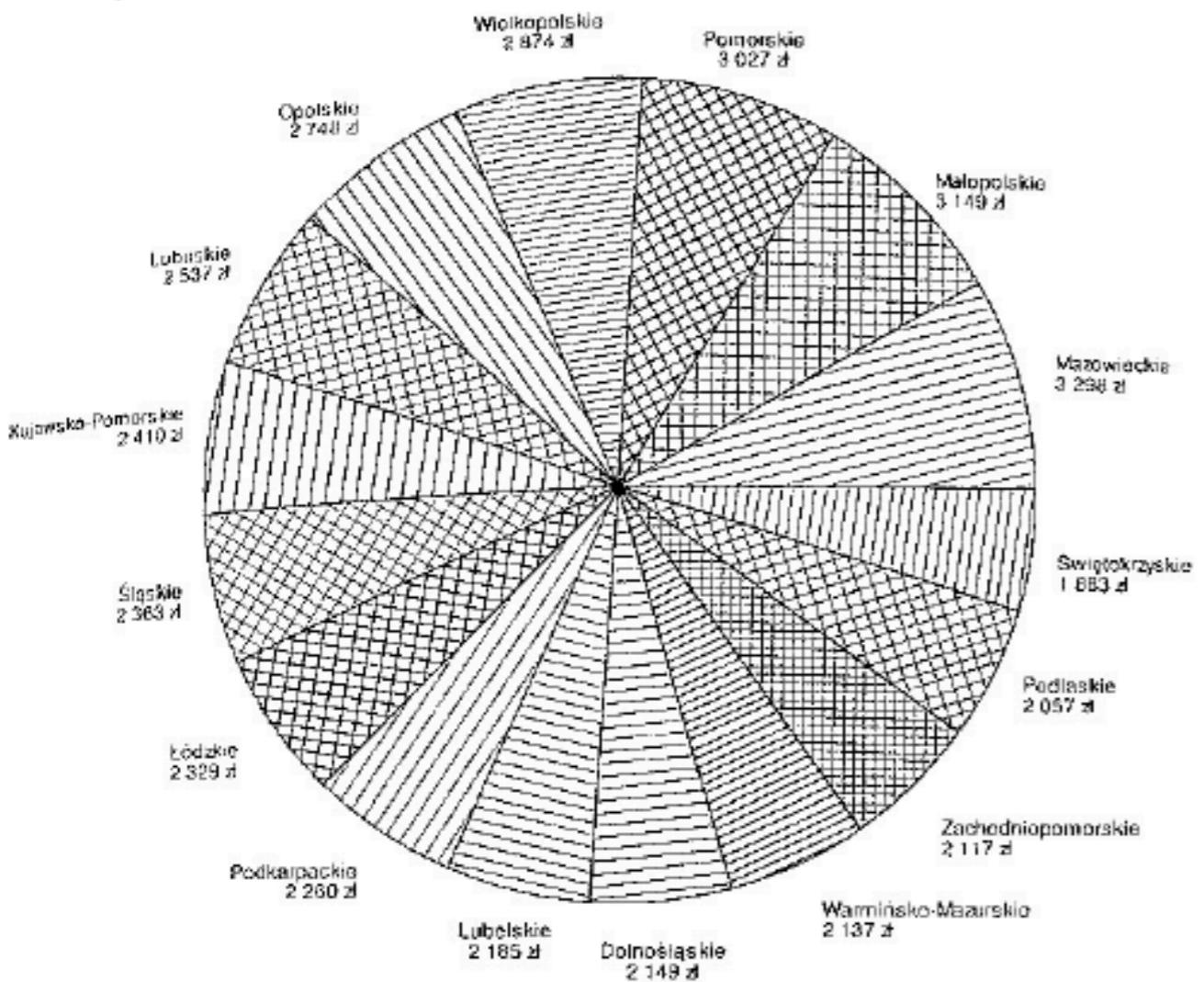
Do not lie

Koszt użytkowania nieruchomości na osobę w gospodarstwie domowym

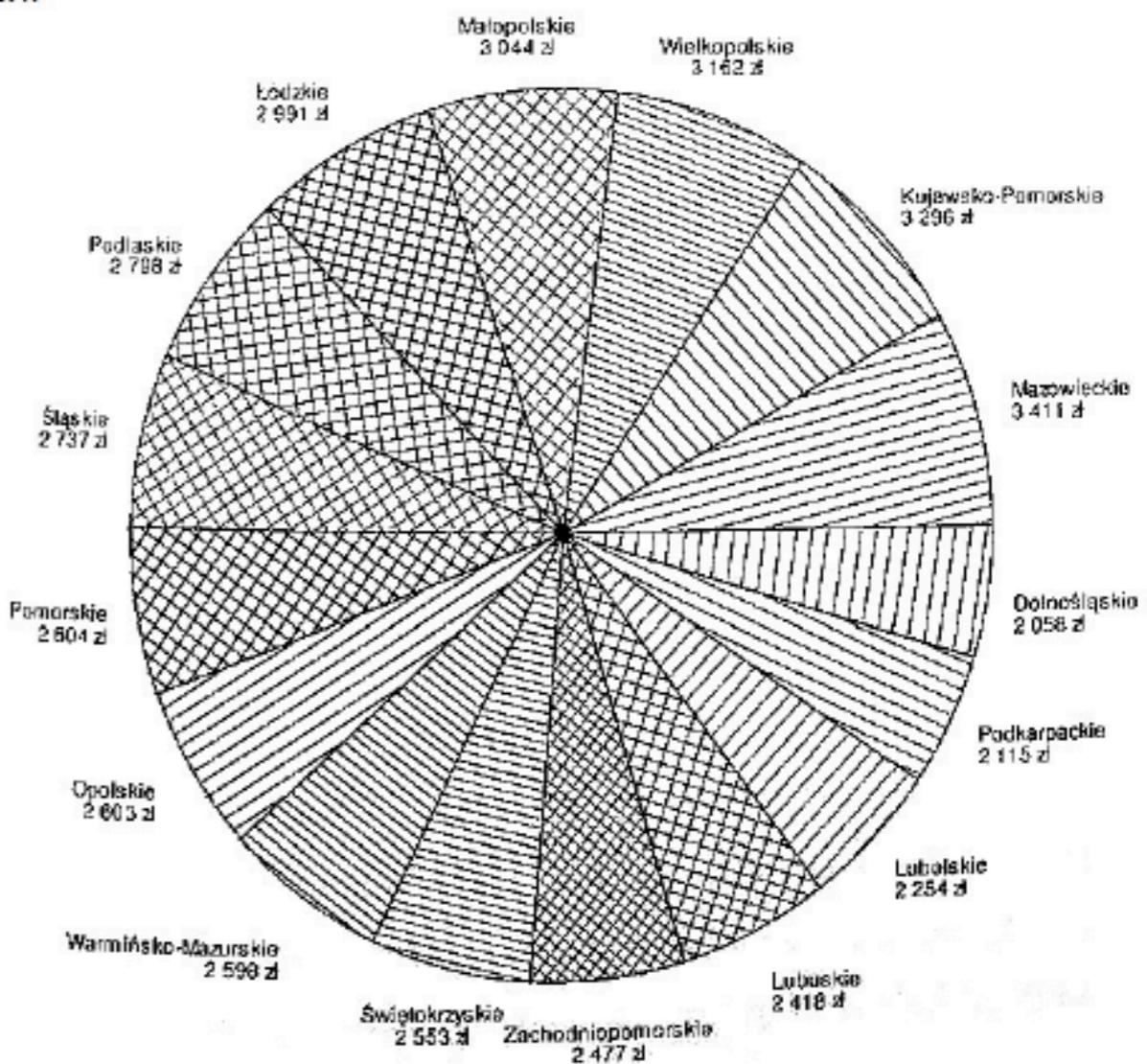


Źródło: Home Broker, GUS

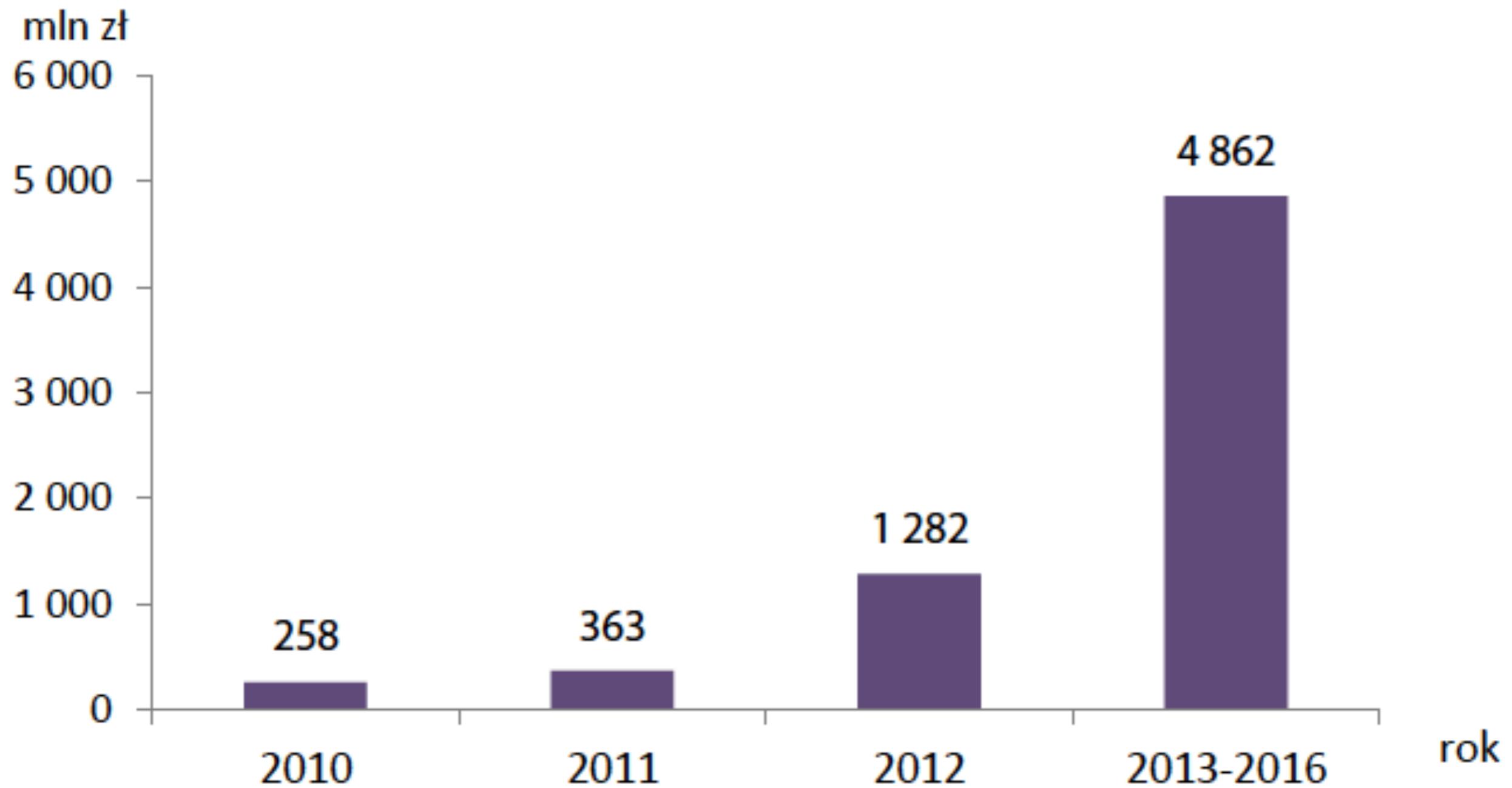
Kobiety



Mężczyźni



I.3 Wydatki deklarowane przez przedsiębiorców na B+R w programach NCBiR w latach 2010-2016.



źródło: Narodowe Centrum Badań i Rozwoju

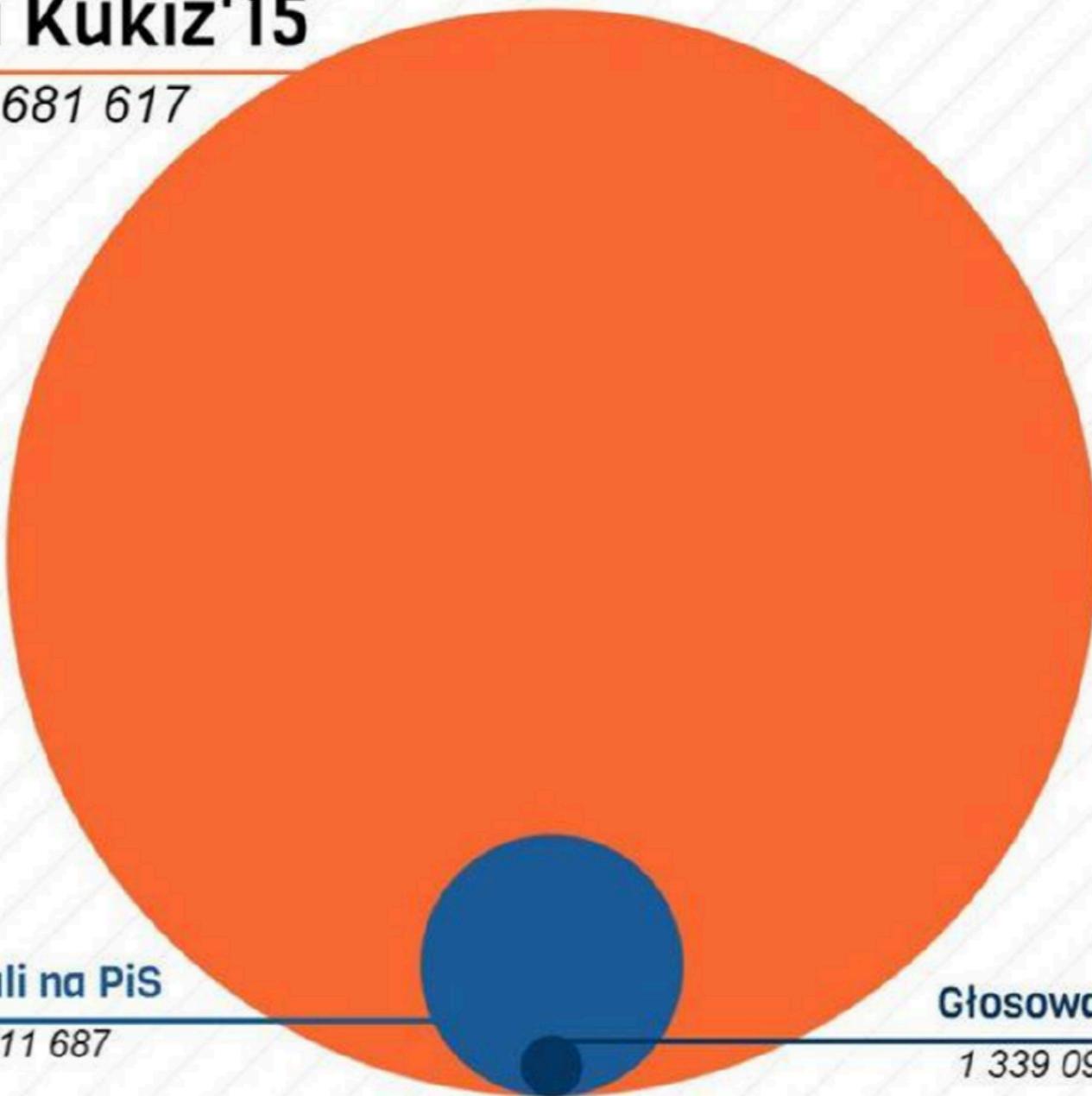
30 732 398 Polaków

UPRAWNIONYCH DO GŁOSOWANIA

* Dane: Państwowa Komisja Wyborcza

**NIE GŁOSOWALI
na PiS i Kukiz'15**

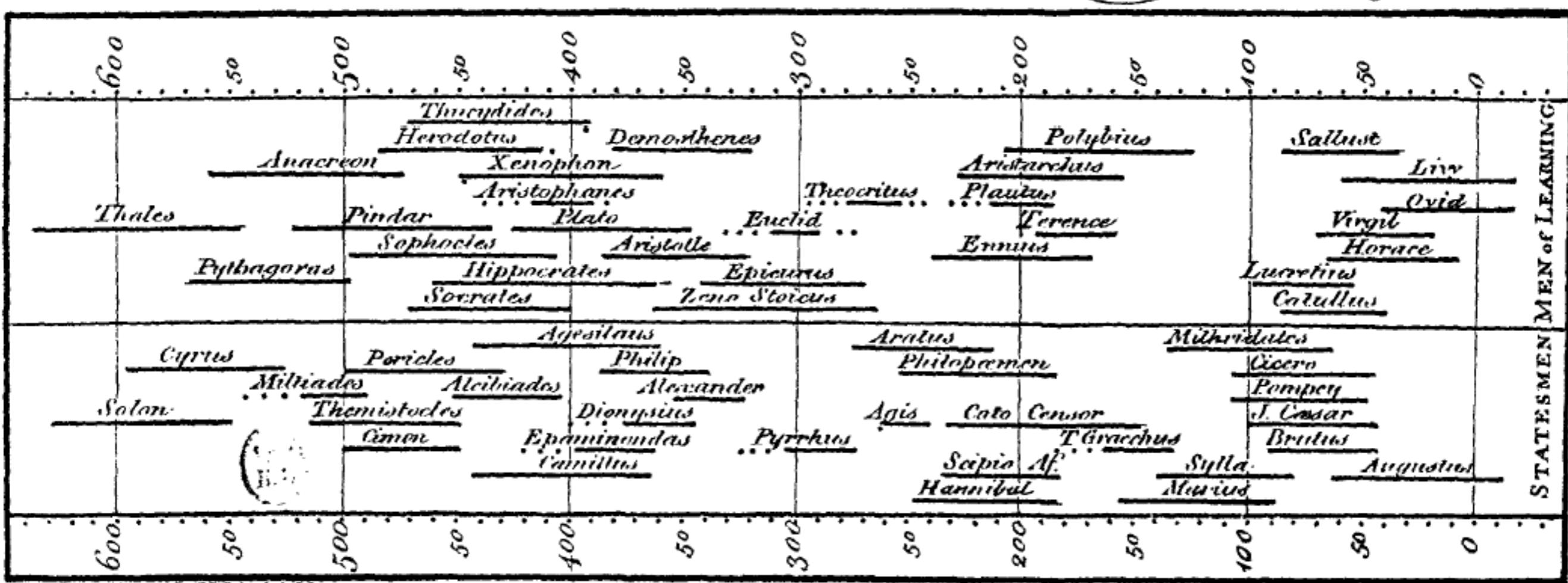
23 681 617



History

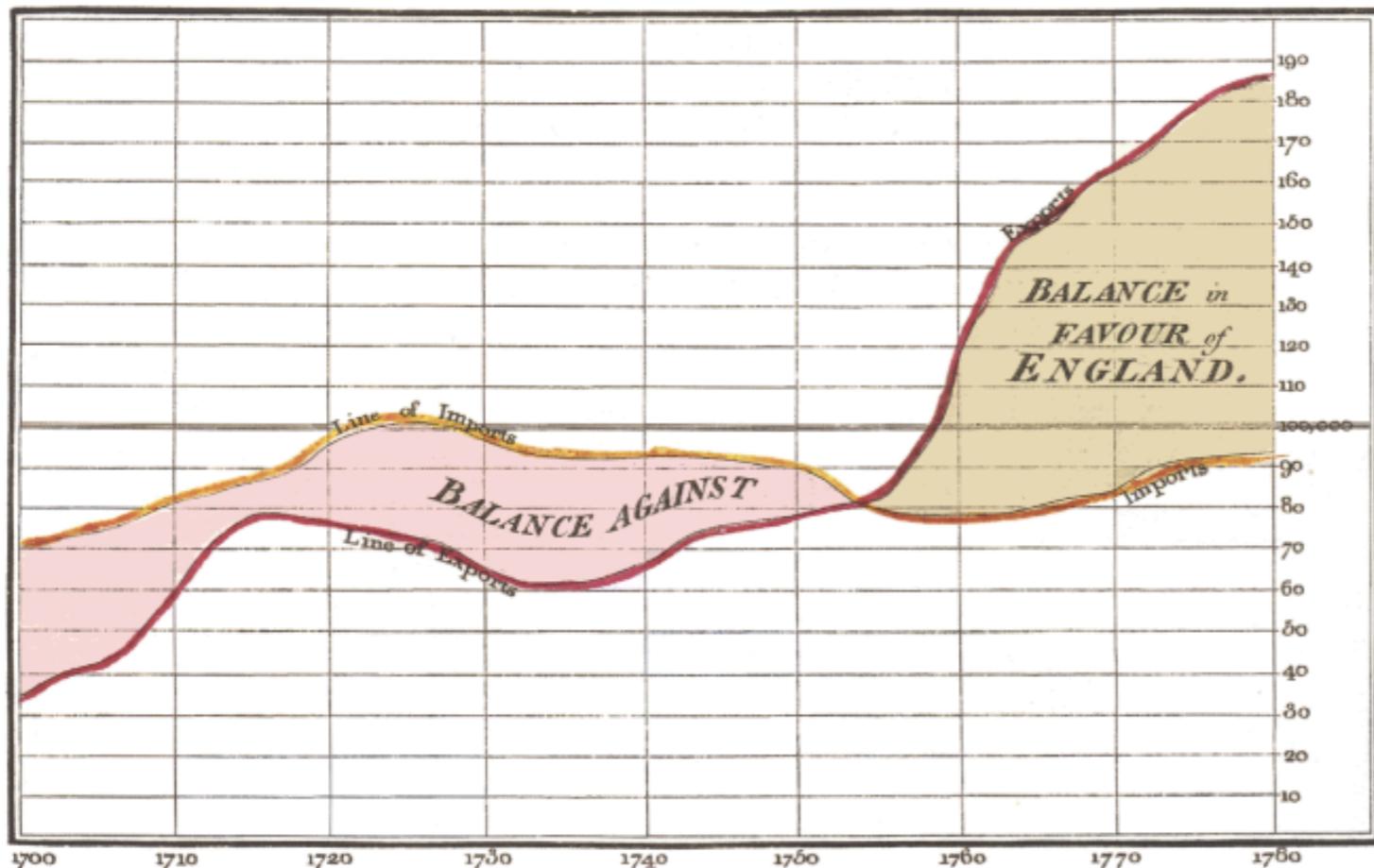
Chart of Biography Joseph Priestley (1765)

A Specimen of a Chart of Biography.



Commercial and Political Atlas William Playfair (1786)

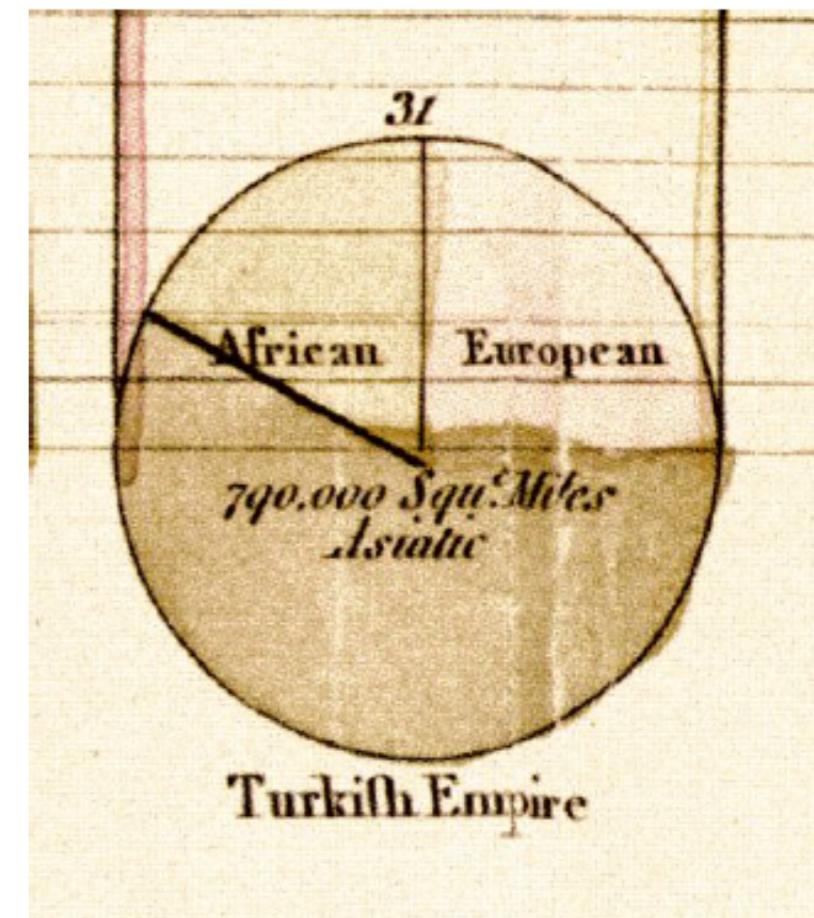
Exports and Imports to and from DENMARK & NORWAY from 1700 to 1780.



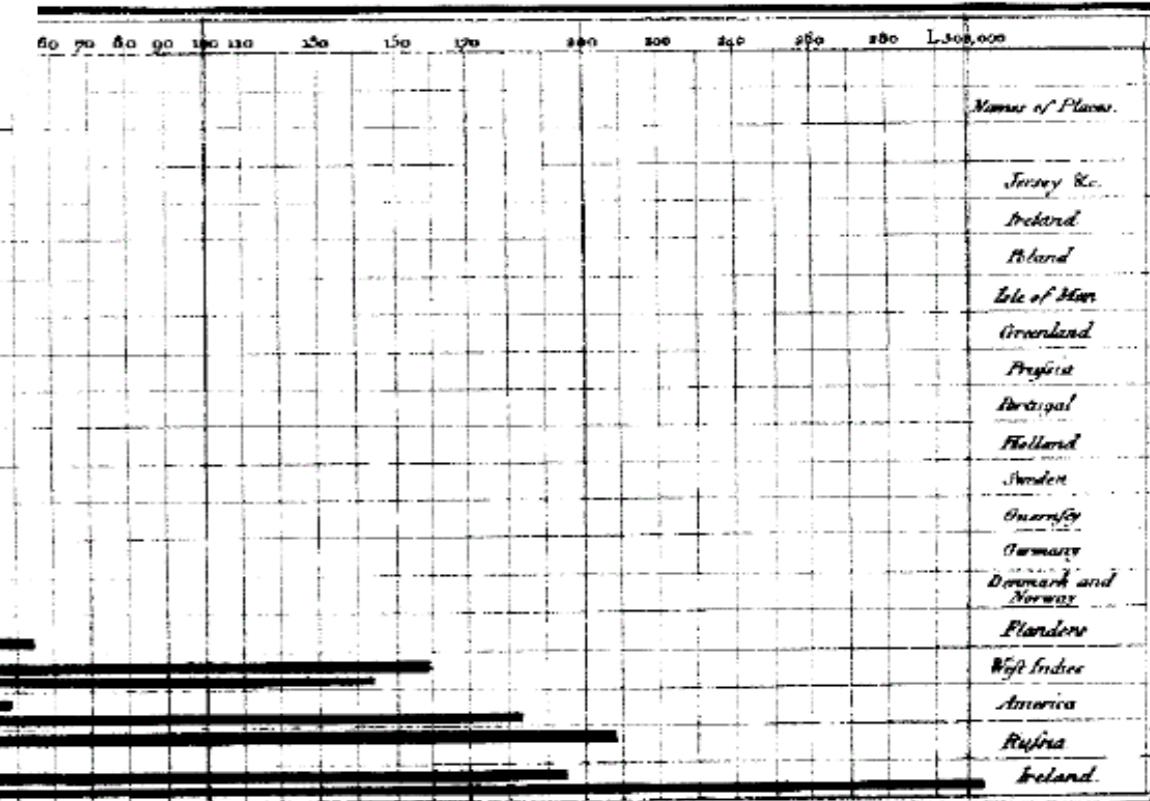
The Bottom line is divided into Years, the Right hand line into £10,000 each.

Published at the Act of Parliament, 10th May 1786, by W^m Playfair

No. 100, Strand, London.



Exports of SCOTLAND to and from different parts for one Year from Christmas 1780 to Christmas 1781.



On the Mode of Communication of Cholera. John Snow. 1855

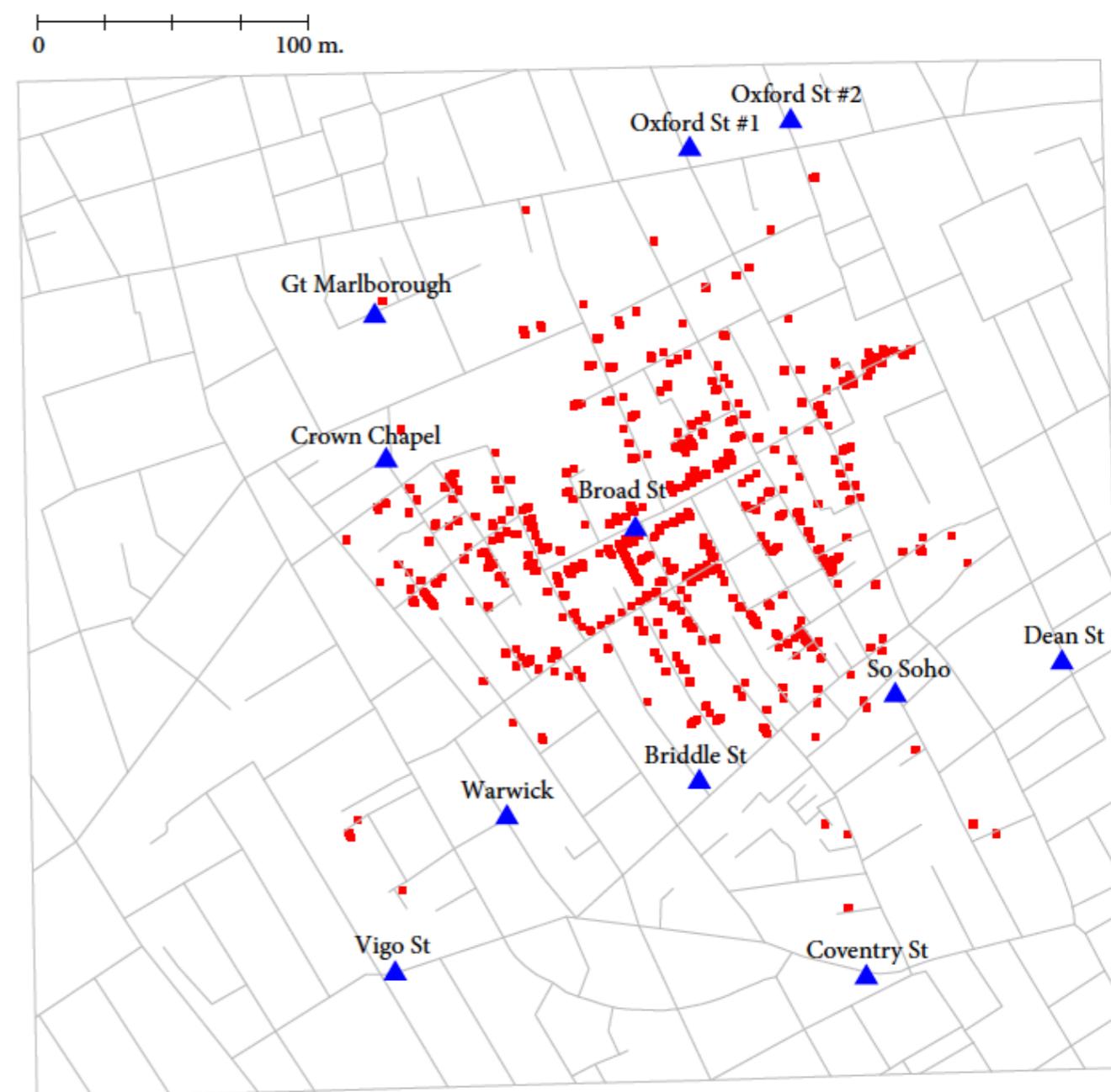
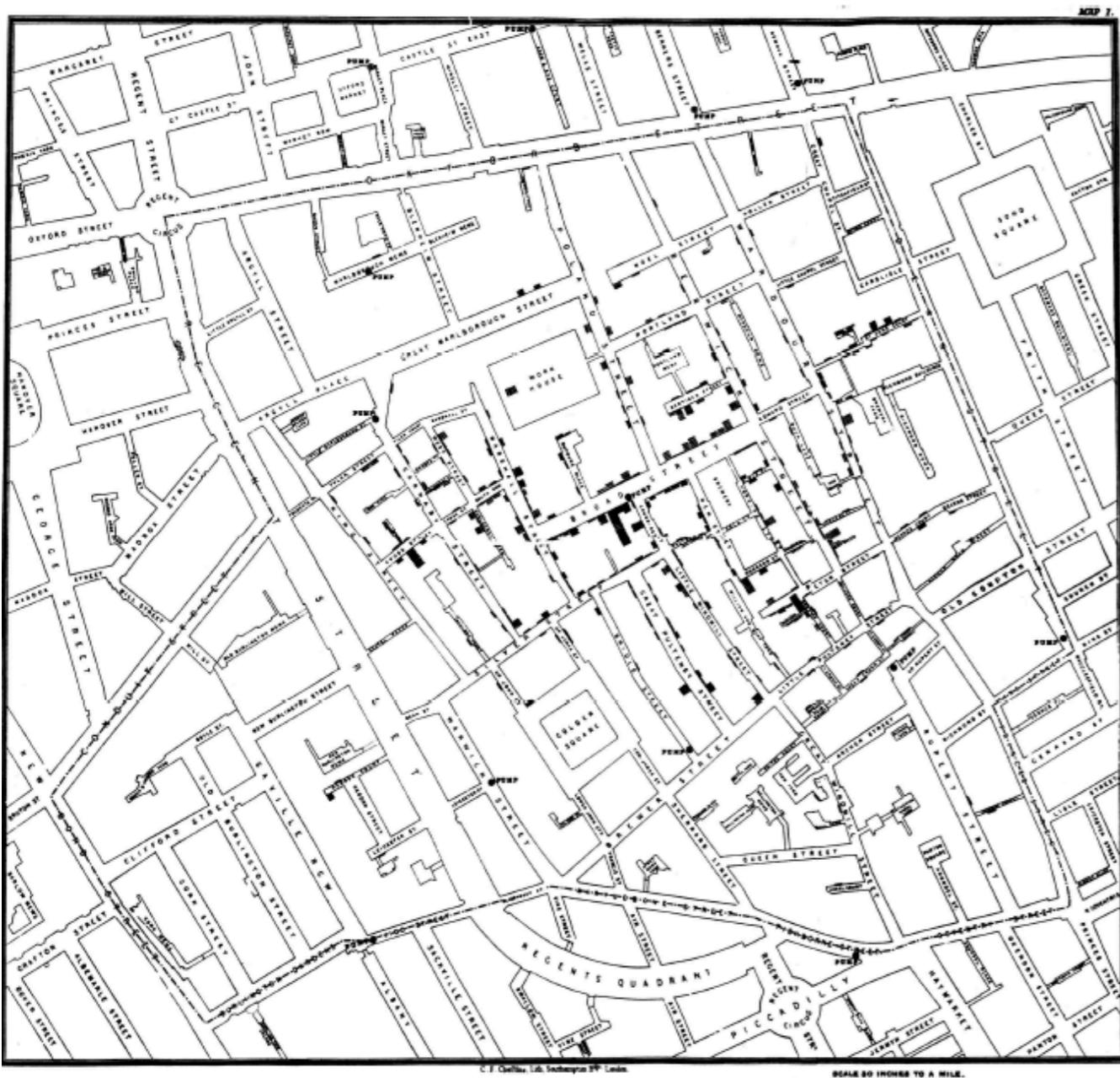
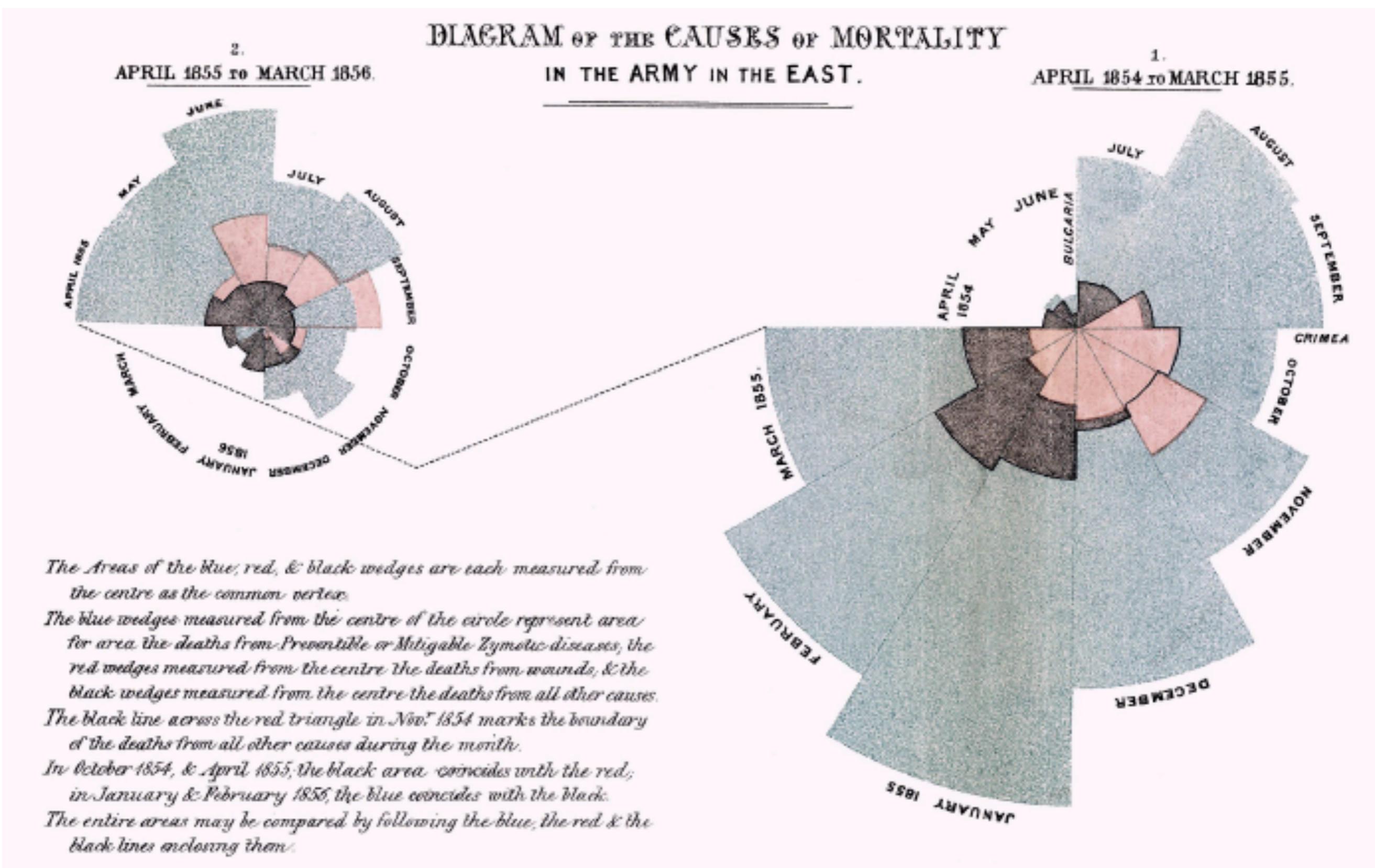


Diagram of the causes of mortality in the army in the East. Florence Nightingale. 1858



More resources



Przemysław Biecek

Odkrywać! Ujawniać! Objaśniać!

Zbiór esejów o sztuce prezentowania danych



Discover! Reveal! Describe!

Essays about the art of data presentation

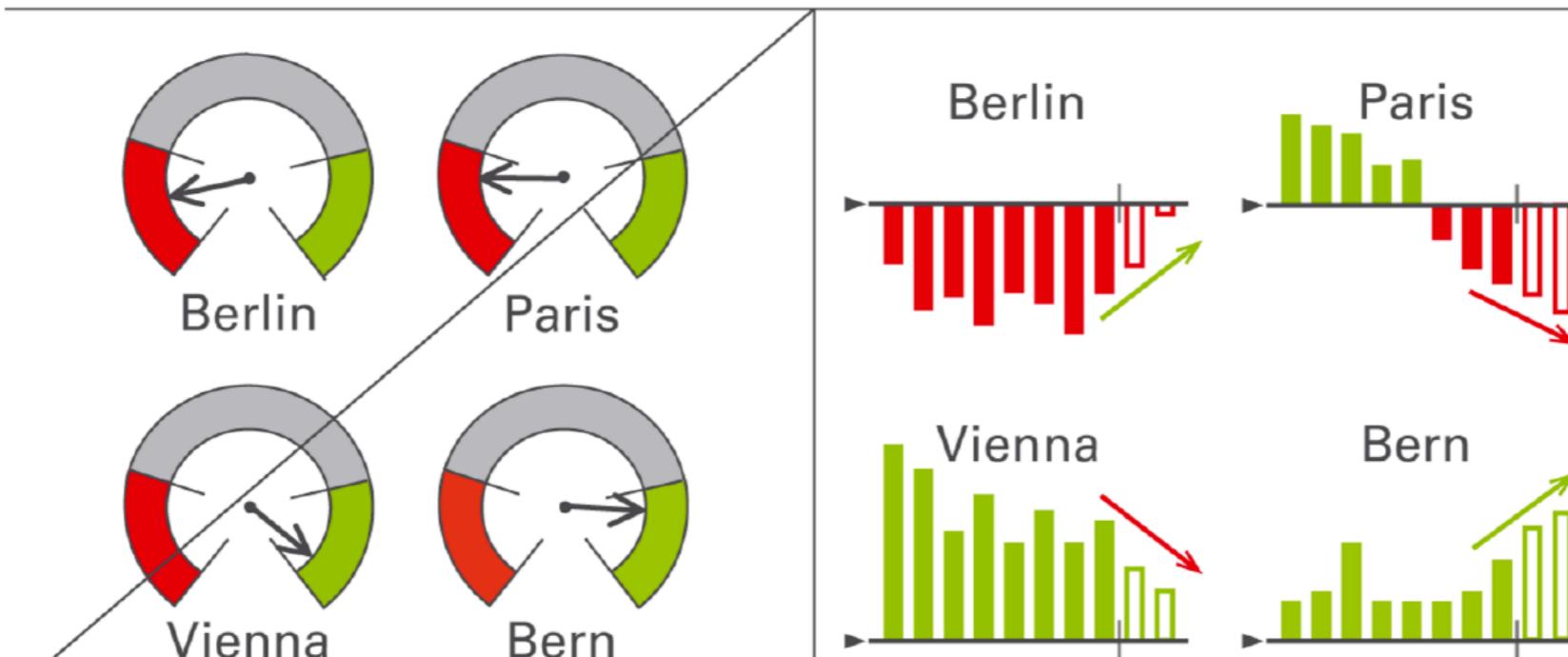
<http://biecek.pl/Eseje>

(english version will be there soon)

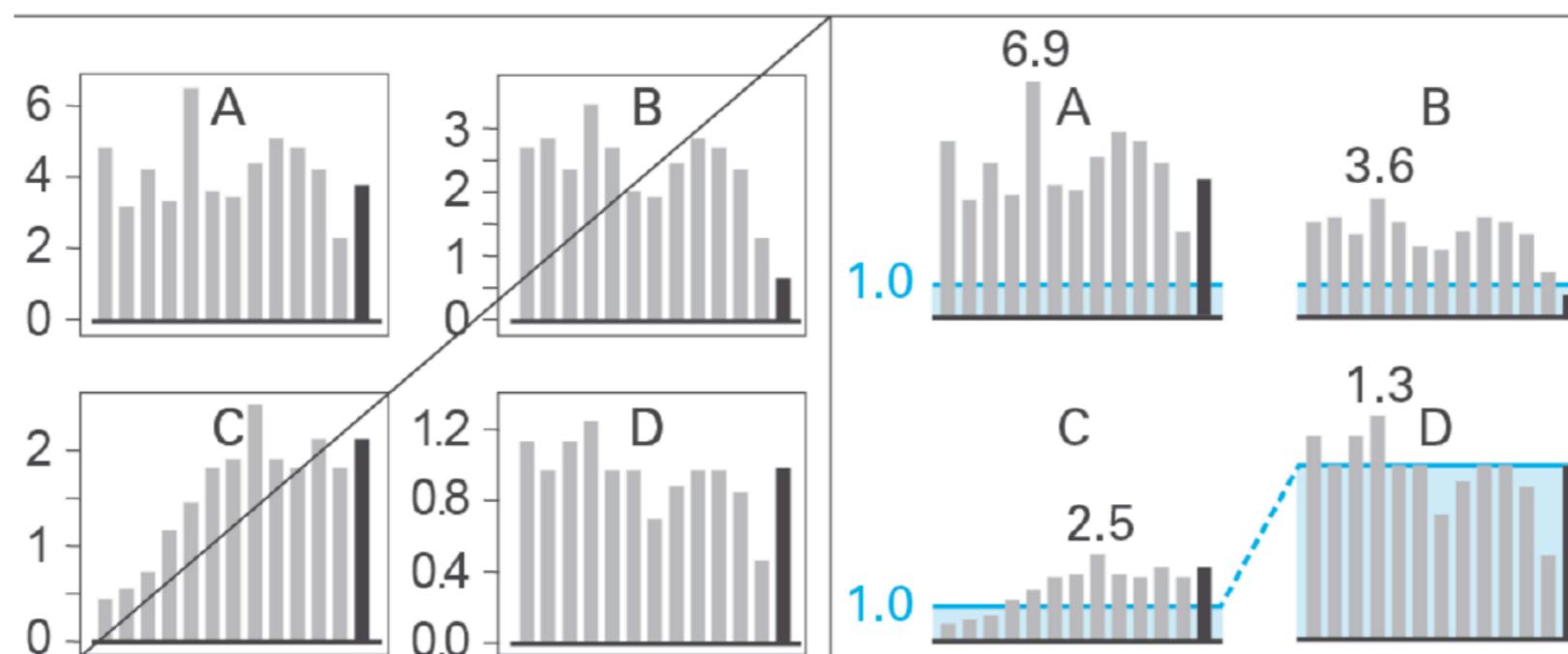
International Business Communication Standards

<http://www.ibcs-a.org/>

EX 2.2 Replace gauges, speedometers



UN 5.2 Unify scaling indicators



Docs ggplot2 <http://docs.ggplot2.org/current/>
 Cookbook for R <http://www.cookbook-r.com/Graphs/>
 Docs ggviz <http://ggviz.rstudio.com/>
 Great blog <http://flowingdata.com/>
 Graphs in NYT <http://kpq.github.io/chartsnthings/>
 Nature Methods, Points of View <http://clearscience.info/wp/?p=546>

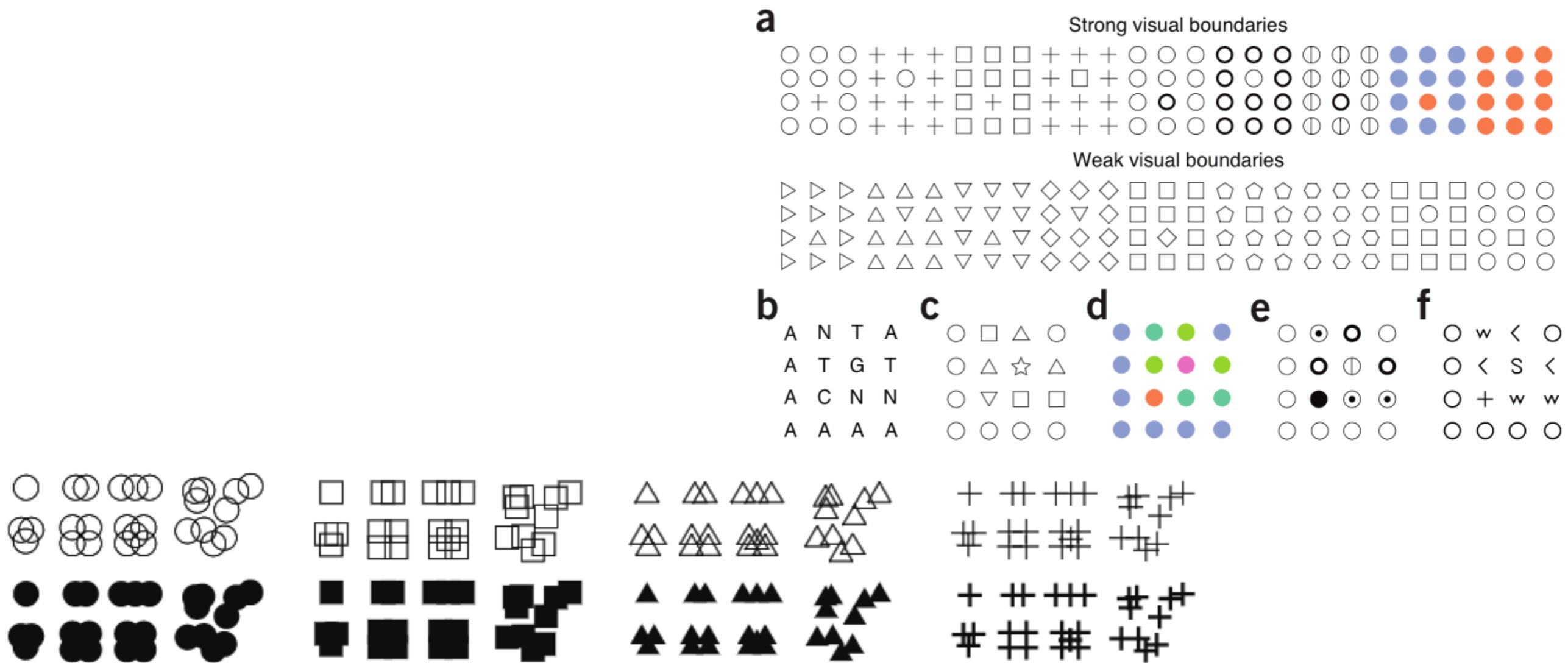


Figure 1 | The hollow circle is a flexible and robust plotting symbol.

[bio project]

MLExpResso

MLExpResso is an R package for integrative analyses and visualisation of gene expression and DNA methylation data.

Key functions of this package are: identification of genes with affected expression, identification of DMR - differentially methylated regions, identification of regions with changes in expression and methylation, visualisation of identified regions.

The package: <https://github.com/geneticsMiNIng/MLGenSig>

The vignette: <https://github.com/geneticsMiNIng/MLGenSig/blob/master/Vignette/Usage.pdf>

<https://github.com/geneticsMiNIng/MLGenSig/blob/master/Vignette/Usage.Rmd>

