# Data Vis

## with ggplot2

## Why and How

http://bit.ly/2xEUGIu

Przemysław Biecek
http://biecek.pl

# Package datasauRus



X Mean: 54.26 32025
Y Mean: 47.83 15781
X SD  : 16.76 50109
Y SD  : 26.93 53144
Corr. : -0.06 45195

https://www.autodeskresearch.com/publications/samestats

# 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)
```
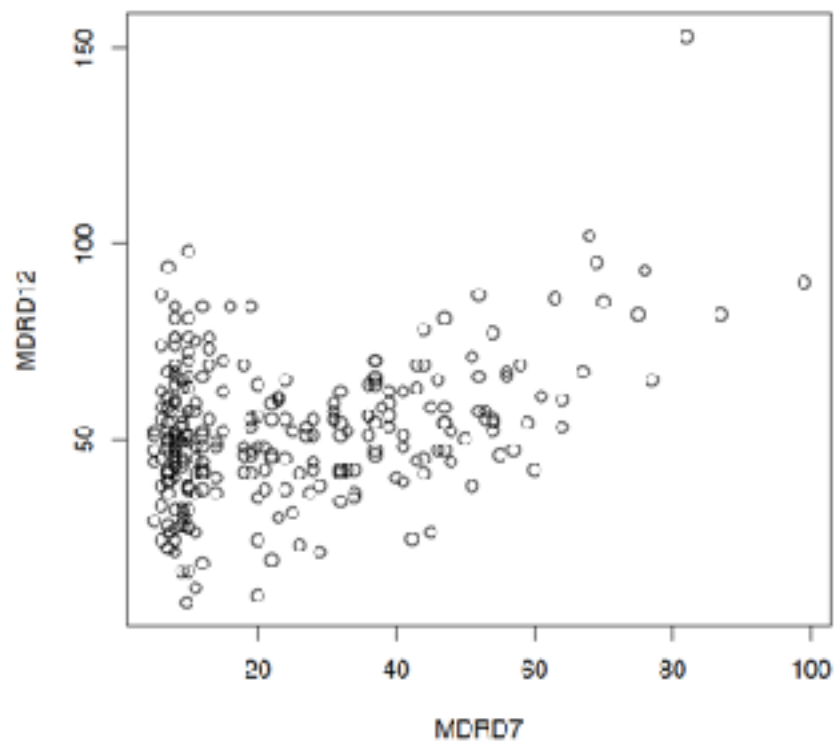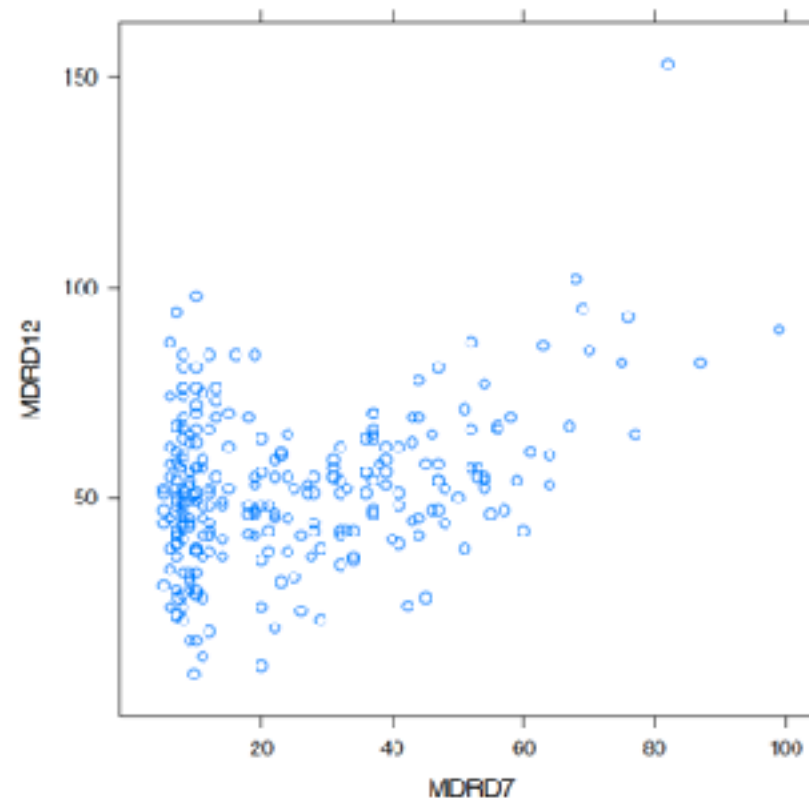
# Find differences between these plots



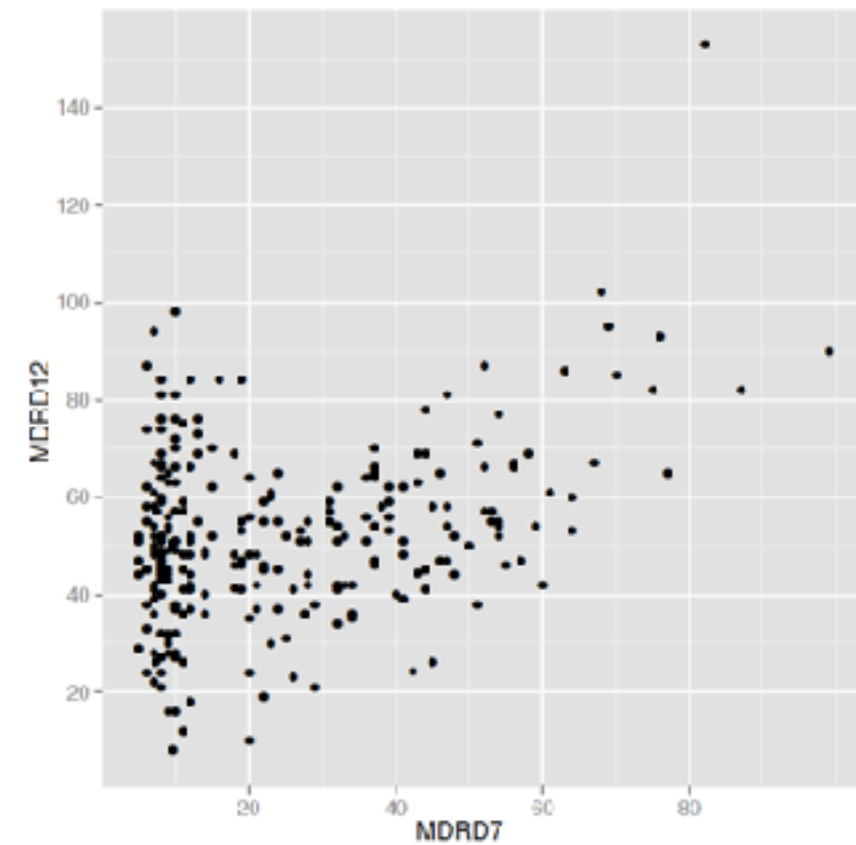**graphics**　　　　　**lattice**　　　　　**ggplot2**
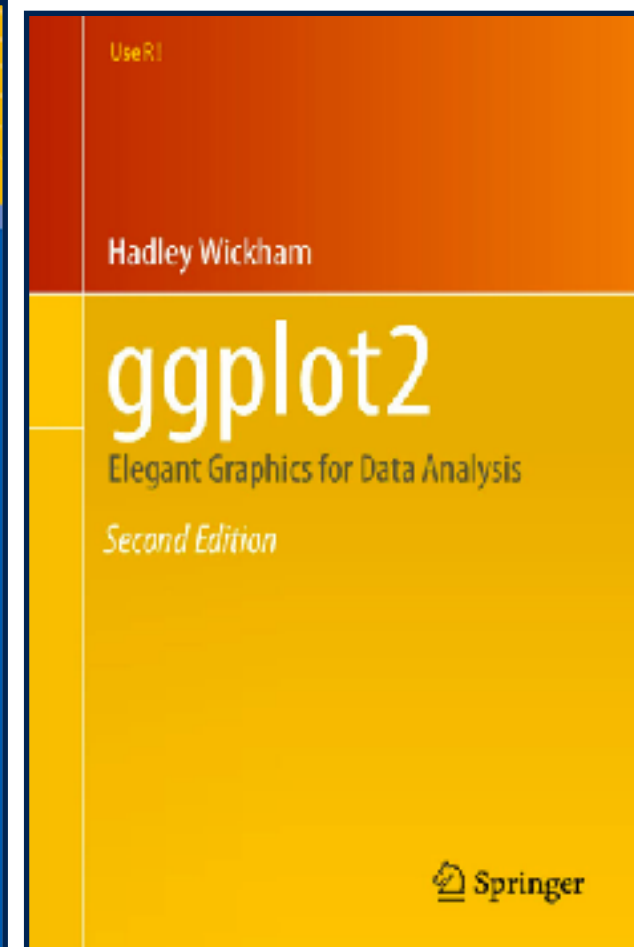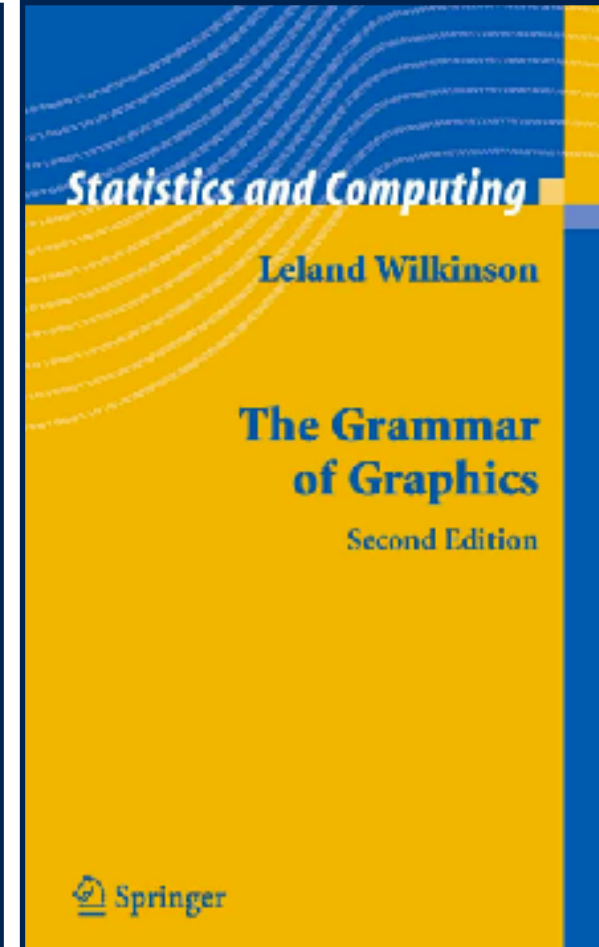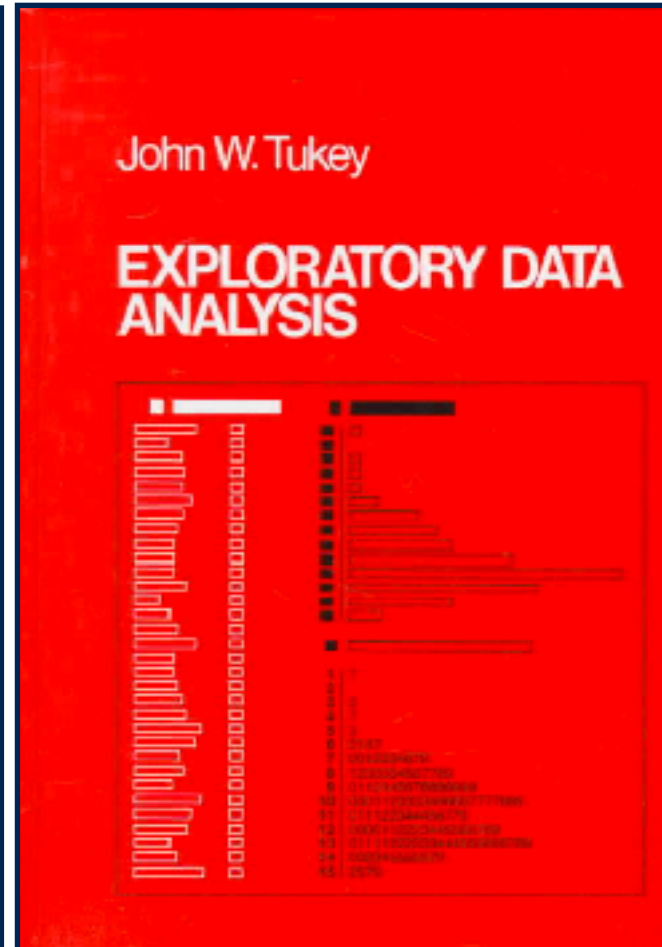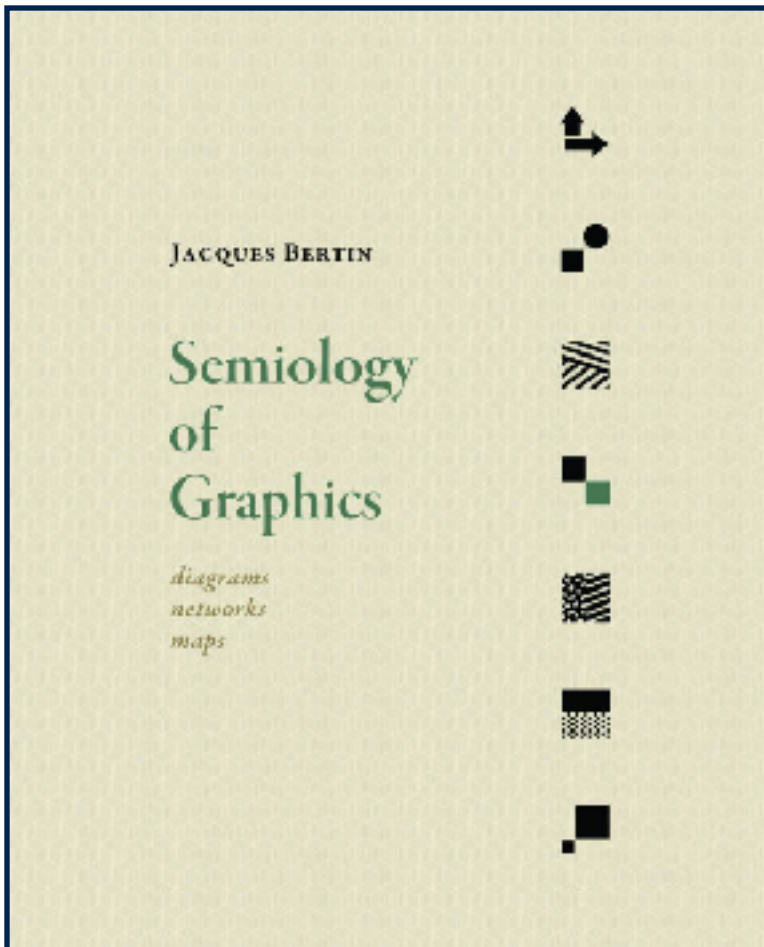
# Why ggplot2?

+ Elegant
+ Highly customisable
+ Uniform
+ Natural
+ Expressive
+ Popular

- Steep learning curve
- Slow
- Evolving pretty fast
           (too fast?)

# Why ggplot2?

**Chart (ggplot)**
Story based on data

**Scales (scale_*)**
Transformation of data variables
to graphical ascetics

*

1

**Skin (theme_*)**
Visual style / theme
of the plot

1

**Coordinate system (coord_*)**
Frames in which the story is based

*

**Panels(facet_*)**
Parts of the story
similar in structure
for data subsets

*

**Layers (layer_*)**
Parts of the story
different structure
same data

1

**Data set (data)**
Table with observations
to be described

1

**Adjustmetns (position_*)**
Position adjustments

*

**Mappings (aes)**
Pairs of data variables and
graphical ascetics

1

**Statsitics (stat_*)**
Data aggregates

1

**Form (geom_*)**
Way in which data is presented
in the layer

http://biecek.pl/Eseje/indexGramatyka.html

# [ hands on live R ]

```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)

serialeIMDB


## (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_ql_4")) %>%
    html_node("table") %>%
    html_table()
  data.frame(Serie = titles[i], tab[,1:4], Season = gsub(tab[,1], pattern="\\..*", replacement=""))
```
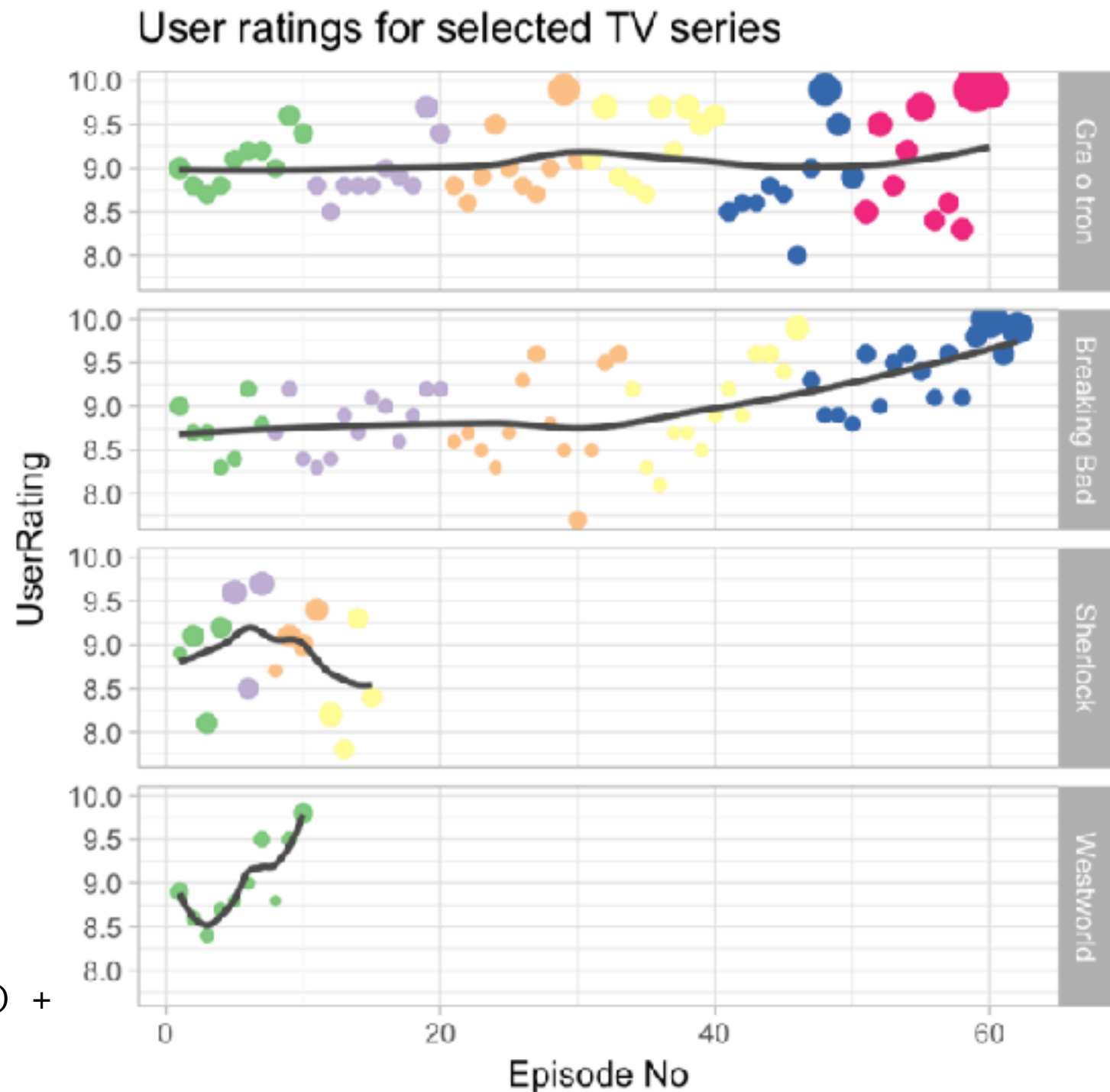
User ratings for selected TV series

```
selected <- c("Gra o tron", "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")
```
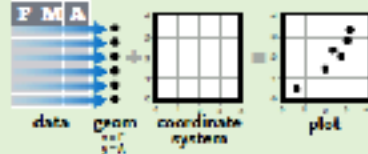
# Data Visualization
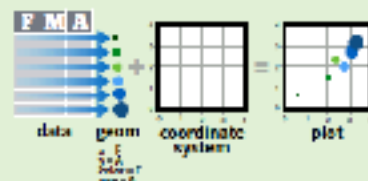## with ggplot2
### Cheat Sheet

**R**Studio

## 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)) +       add layers,
  geom_point(aes(color = cyl)) +   elements with +
  geom_smooth(method ="lm") +      layer = geom +
  coord_cartesian() +              default stat +
  scale_color_gradient() +         layer specific
  theme_bw()                       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")

 **a + geom_density(kernel = "gaussian")**
x, y, alpha, color, fill, linetype, size, weight
b + geom_density(aes(y = ..county..))

 **a + geom_dotplot()**
x, y, alpha, color, fill

 **a + geom_freqpoly()**
x, y, alpha, color, linetype, size
b + geom_freqpoly(aes(y = ..density..))

 **a + geom_histogram(binwidth = 5)**
x, y, alpha, color, fill, linetype, size, weight
b + geom_histogram(aes(y = ..density..))

#### 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

#### Discrete X, Continuous Y
g <- ggplot(mpg, aes(class, hwy))

 **g + geom_bar(stat = "identity")**
x, y, alpha, color, fill, linetype, size, weight

 **g + geom_boxplot()**
lower, middle, upper, x, ymax, ymin, alpha, color, fill, linetype, shape, size, weight

 **g + geom_dotplot(binaxis = "y",**
**stackdir = "center")**
x, y, alpha, color, fill

 **g + geom_violin(scale = "area")**
x, y, alpha, color, fill, linetype, size, weight

#### Discrete X, Discrete Y
h <- ggplot(diamonds, aes(cut, color))

 **h + geom_jitter()**
x, y, alpha, color, fill, shape, size

#### 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 <- ggplot(df, aes(grp, fit, ymin = fit-se, ymax = fit+se))

 **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 <- ggplot(data, aes(fill = murder))

 **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_contour(aes(z = z))**
x, y, z, alpha, colour, linetype, size, weight

 **m + geom_raster(aes(fill = z), hjust=0.5,**
**vjust=0.5, interpolate=FALSE)**
x, y, alpha, fill

 **m + geom_tile(aes(fill = z))**
x, y, alpha, color, fill, linetype, size

Learn more at docs.ggplot2.org · ggplot2 0.9.3.1 · Updated 3/15

https://www.rstudio.com/wp-content/uploads/2015/03/ggplot2-cheatsheet.pdf

# Do not lie

# Koszt użytkowania nieruchomości na osobę w gospodarstwie domowym



| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 205 zł | 206 zł | 206 zł | 207 zł | 207 zł | 208 zł | 209 zł | 210 zł | 211 zł | 212 zł |
| 2011 I | 2011 II | 2011 III | 2011 IV | 2011 V | 2011 VI | 2011 VII | 2011 VIII | 2011 IX | 2011 X |

Źródło: Home Broker, GUS

Kobiety

Wielkopolskie 2 874 zł
Pomorskie 3 027 zł
Opolskie 2 748 zł
Małopolskie 3 149 zł
Lubuskie 2 537 zł
Mazowieckie 3 298 zł
Kujawsko-Pomorskie 2 410 zł
Świętokrzyskie 1 883 zł
Śląskie 2 363 zł
Podlaskie 2 057 zł
Łódzkie 2 329 zł
Zachodniopomorskie 2 117 zł
Podkarpackie 2 260 zł
Warmińsko-Mazurskie 2 137 zł
Lubelskie 2 185 zł
Dolnośląskie 2 149 zł

Mężczyźni

Małopolskie 3 044 zł
Wielkopolskie 3 162 zł
Łódzkie 2 991 zł
Kujawsko-Pomorskie 3 296 zł
Podlaskie 2 798 zł
Mazowieckie 3 411 zł
Śląskie 2 737 zł
Dolnośląskie 2 058 zł
Pomorskie 2 604 zł
Podkarpackie 2 115 zł
Opolskie 2 603 zł
Lubelskie 2 254 zł
Warmińsko-Mazurskie 2 598 zł
Lubuskie 2 418 zł
Świętokrzyskie 2 553 zł
Zachodniopomorskie 2 477 zł

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

**NIE GŁOSOWALI na PiS i Kukiz'15**
*23 681 617*

**Głosowali na PiS**
5 711 687

**Głosowali na Kukiz'15**
1 339 094

# History

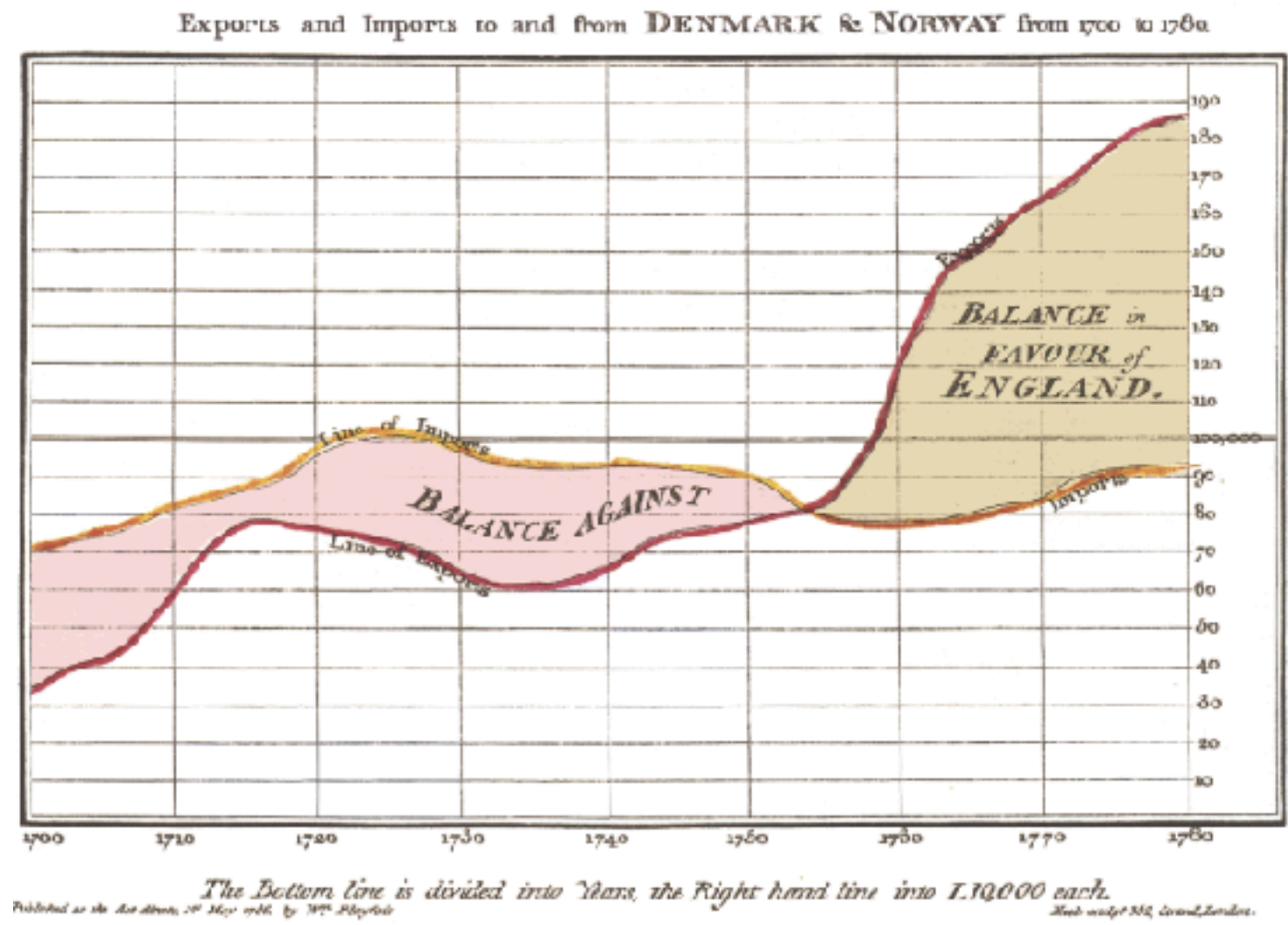A Specimen of a Chart of Biography.

MEN of LEARNING

Thucydides
Herodotus · Demosthenes
Anacreon · Xenophon · Polybius · Sallust
Aristophanes · Aristarchus · Livy
Thales · Pindar · Plato · Theocritus · Plautus · Ovid
Sophocles · Aristotle · Euclid · Terence · Virgil
Pythagoras · Hippocrates · Epicurus · Ennius · Horace
Socrates · Zeno Stoicus · Lucretius
Catullus

STATESMEN

Agesilaus · Aratus · Mithridates
Cyrus · Pericles · Philip · Philopæmen · Cicero
Miltiades · Alcibiades · Alexander · Pompey
Solon · Themistocles · Dionysius · Agis · Cato Censor · J. Cæsar
Cimon · Epaminondas · Pyrrhus · T. Gracchus · Brutus
Camillus · Scipio Af. · Sylla · Augustus
Hannibal · Marius

600 · 50 · 500 · 50 · 400 · 50 · 300 · 50 · 200 · 50 · 100 · 50 · 0

J. Priestley LL.D F.R.S. inv.t et del.

# Commercial and Political Atlas William Playfair (1786)



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

BALANCE in FAVOUR of ENGLAND.

Line of Imports

BALANCE AGAINST

Line of Exports

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



31

African   European

790,000 Squ. Miles Asiatic

Turkish Empire

# 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

Discover! Reveal! Describe!

Essays about the art of data presentation
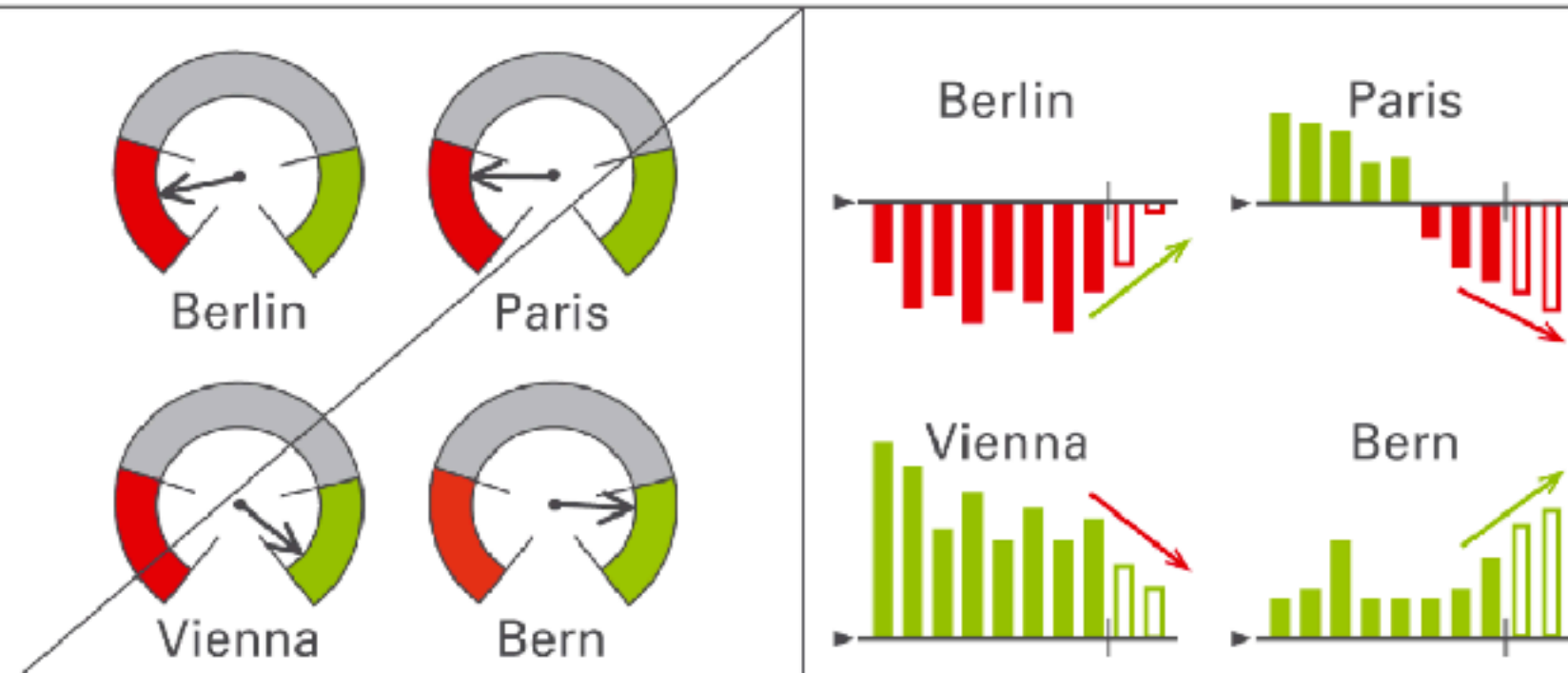
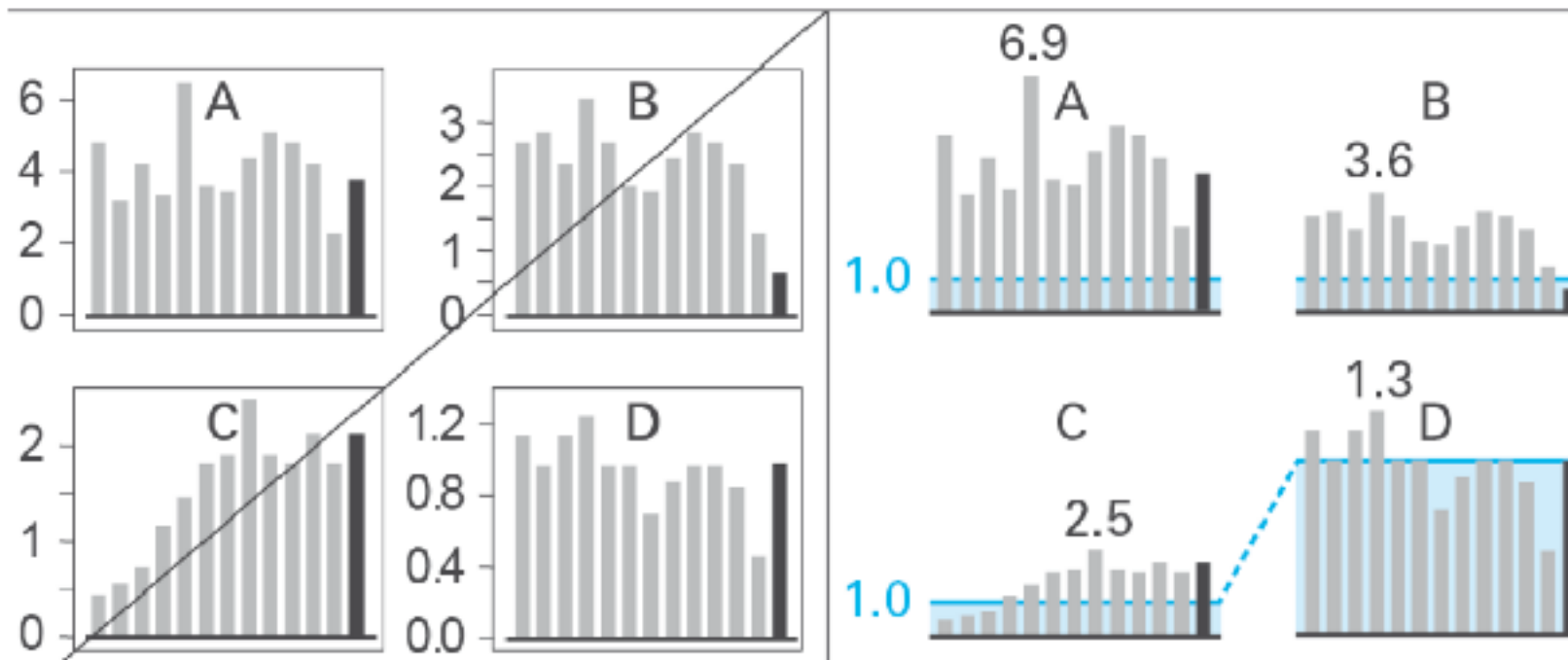http://biecek.pl/Eseje

(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 ggvis http://ggvis.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 ]