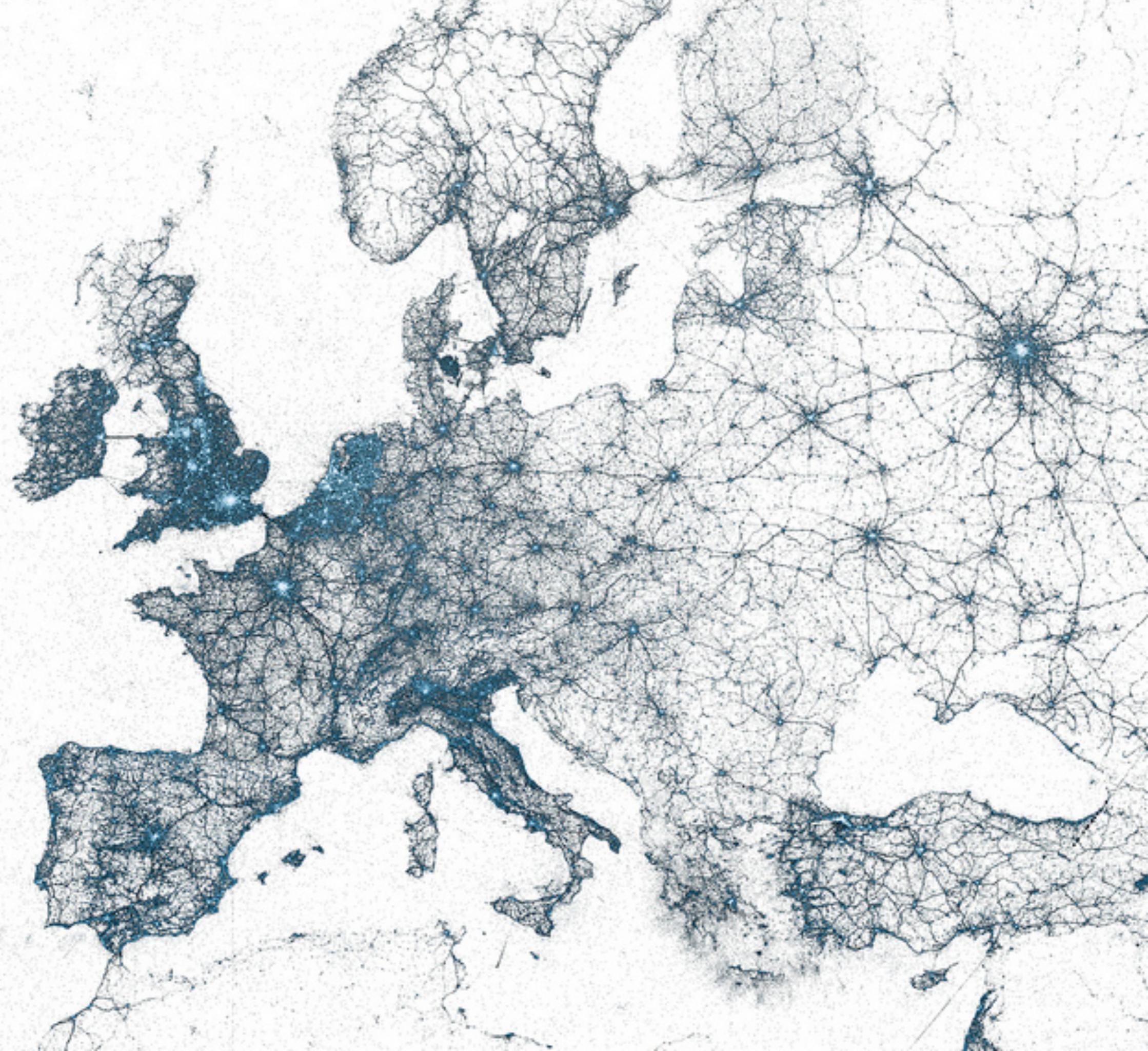


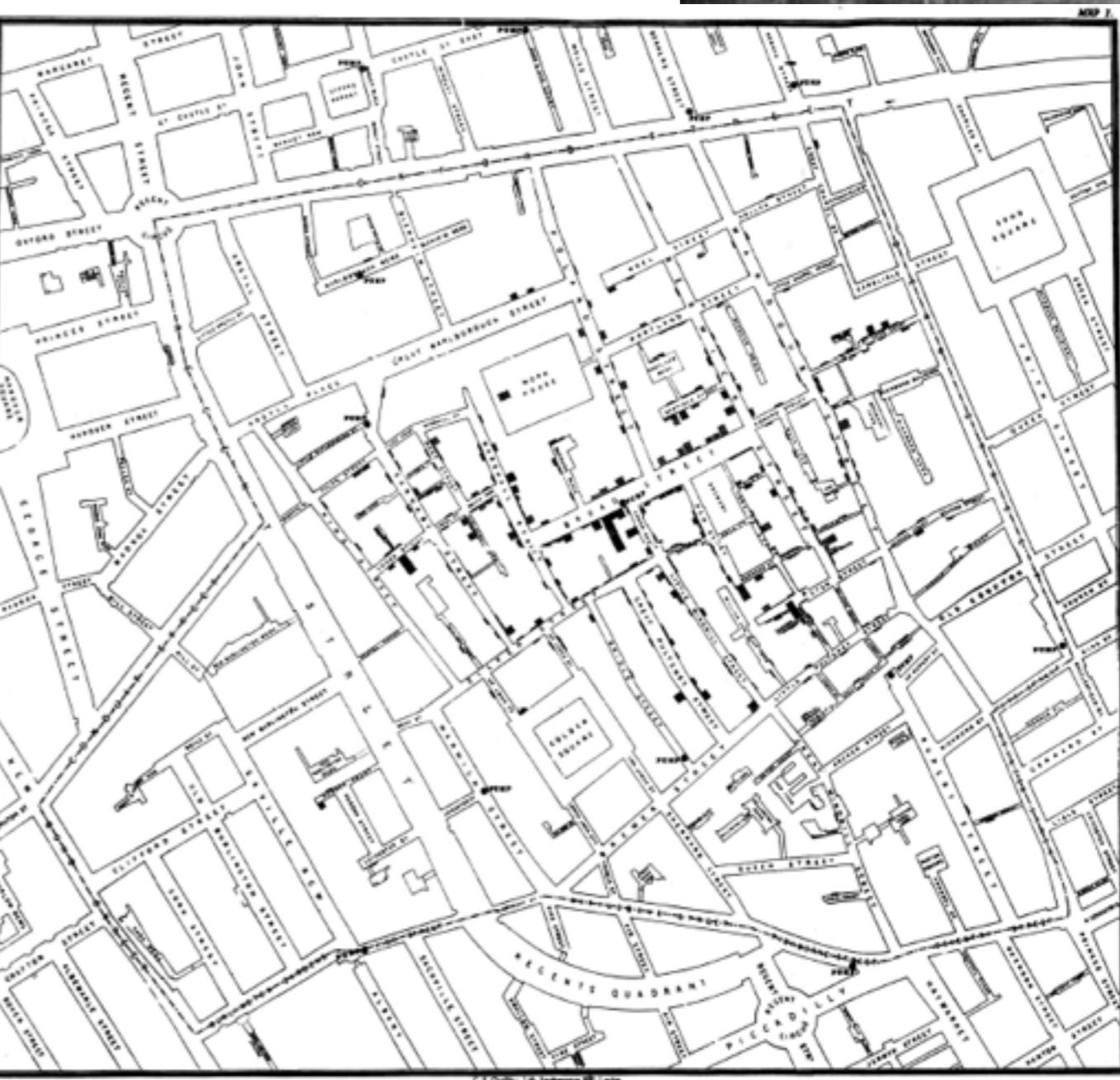
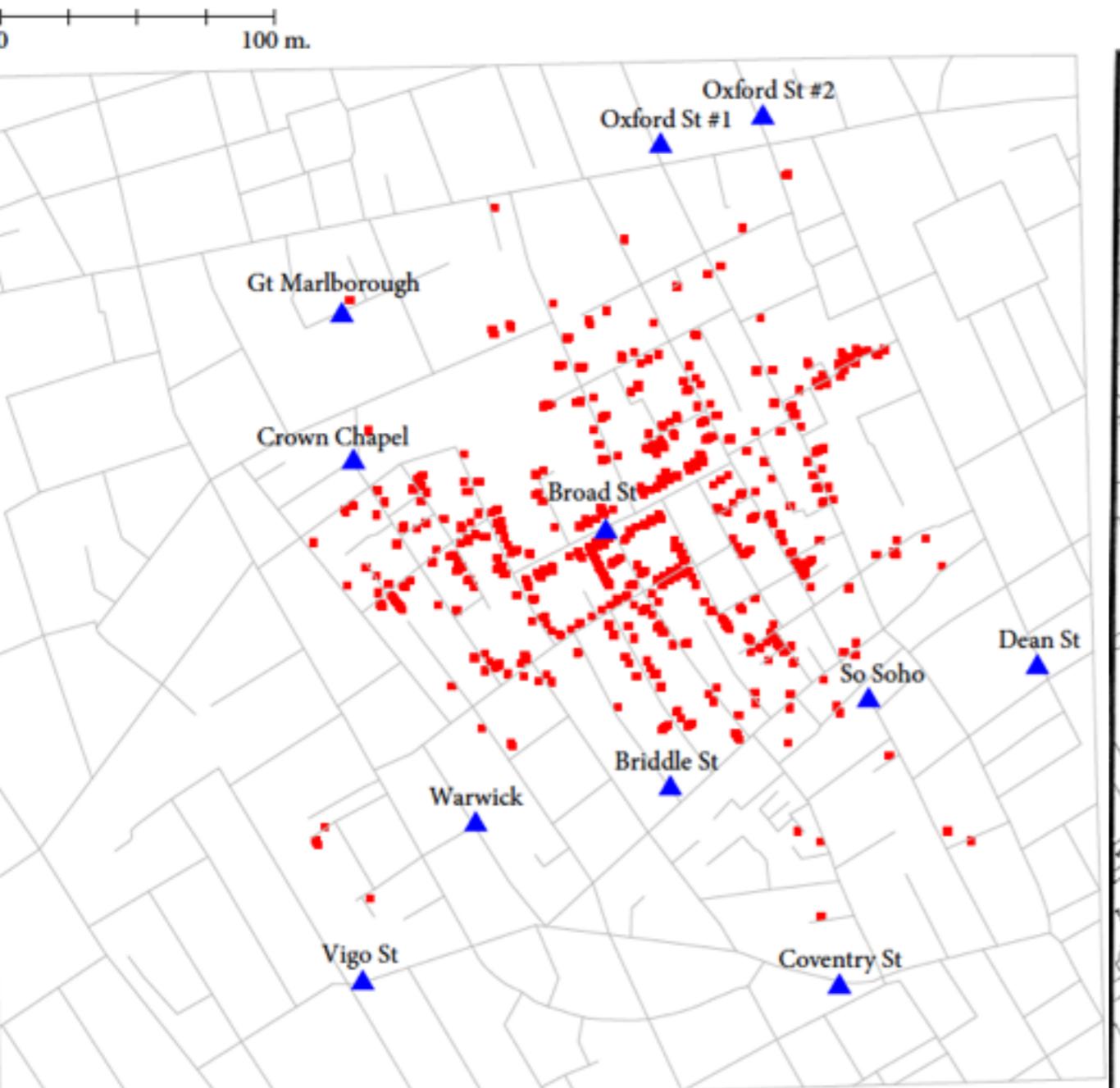
Wizualizacja danych. Co robić, czego nie robić, dlaczego warto poznać pakiet ggplot2 dla R

Przemysław Biecek
Uniwersytet Warszawski
Politechnika Warszawska

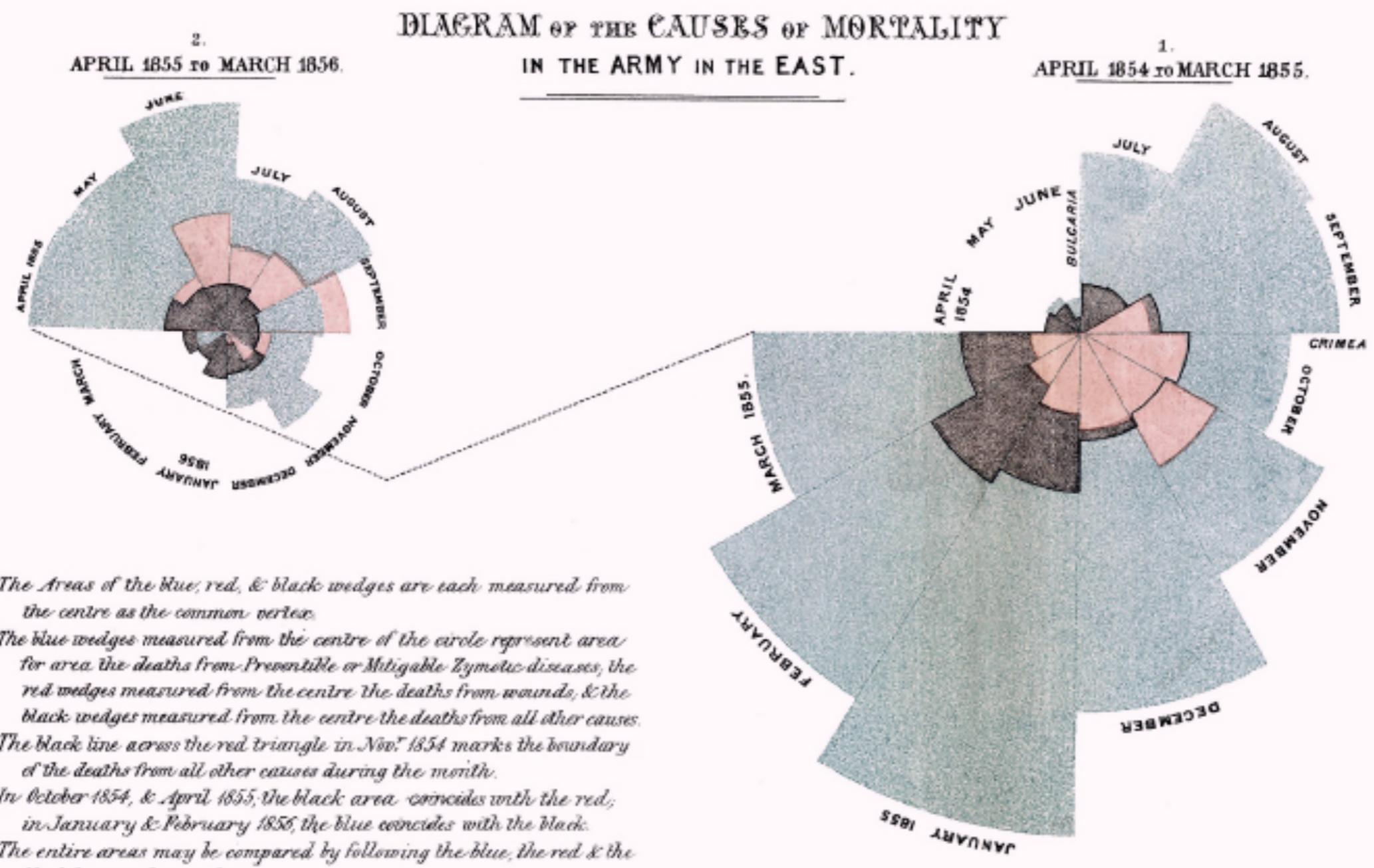
Czy wizualizacja danych jest przydatna?



Kim jest John Snow?

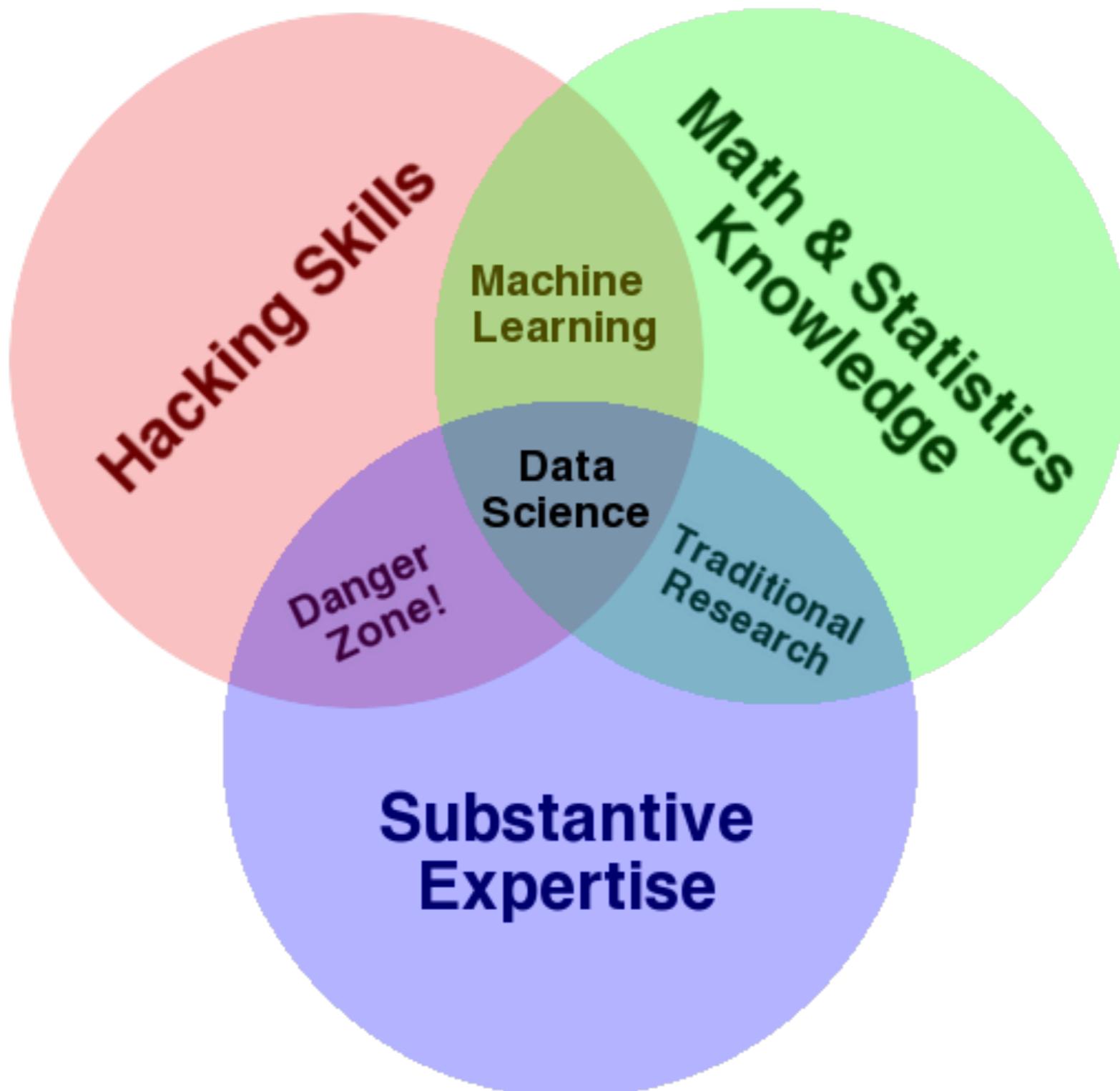


Kim jest Florence Nightingale?



Ale czy to jest Data Science?

Zagadka 1: Gdzie jest wizualizacja danych?



Zagadka 2: Tabela? Wykres? Słownie?

Liczba ludności Polski	W tym kobiet	W tym mężczyzn
38 511 800	51,6%	48,4%

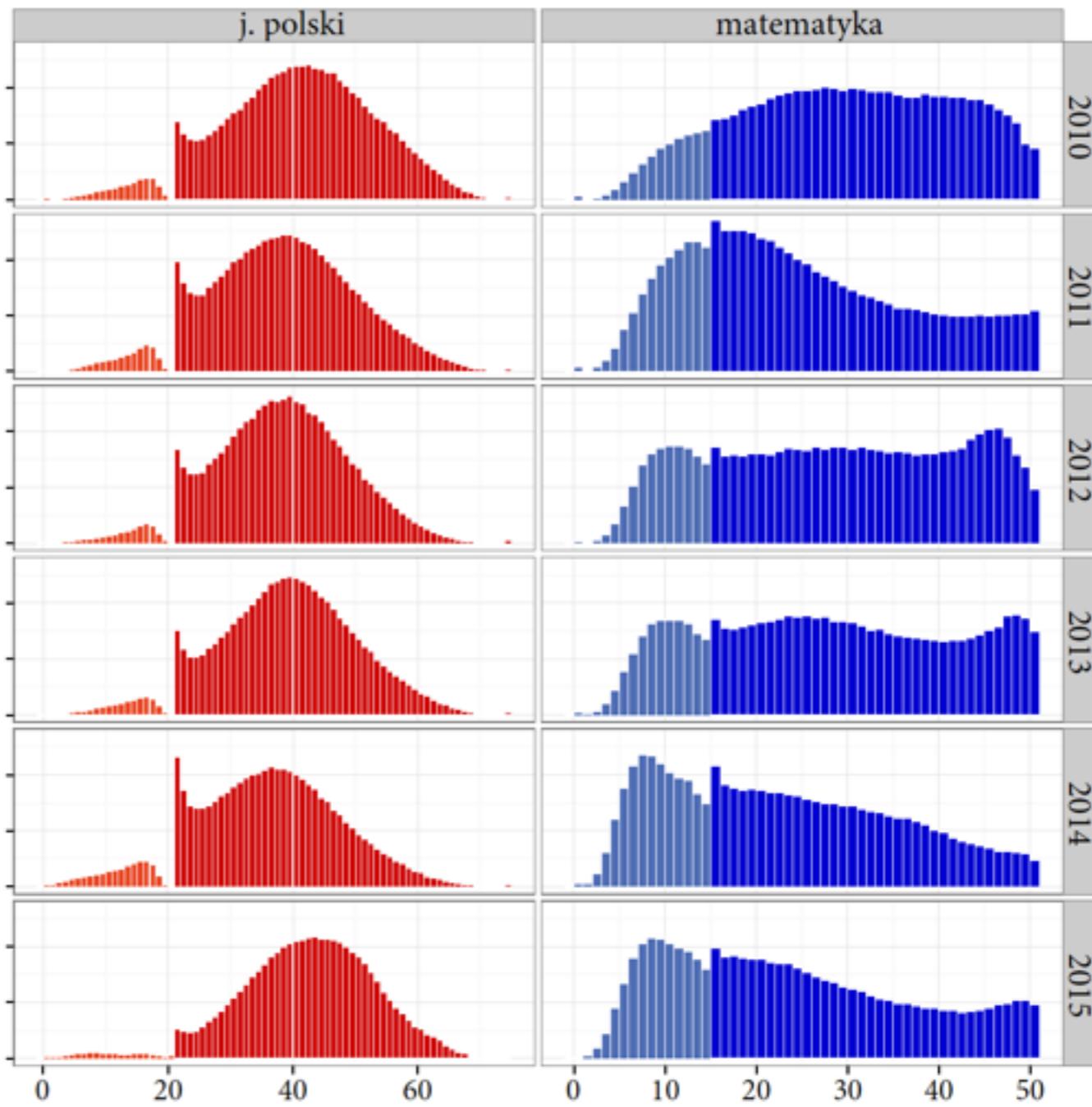


W wyniku przeprowadzenia Narodowego Spisu Powszechnego w roku 2011 ustalono, że w Polsce mieszka 38 511 800 osób, z czego 48,4% to mężczyźni, a 51,6% to kobiety.

Zagadka 2: Tabela? Wykres? Słownie?

punkty	przedmiot	2010	2011	2012	2013	2014	2015
...
6	j. polski	0,09	0,09	0,09	0,09	0,25	0,16
7	j. polski	0,12	0,14	0,11	0,12	0,28	0,16
8	j. polski	0,16	0,18	0,12	0,14	0,34	0,19
9	j. polski	0,19	0,22	0,14	0,19	0,36	0,19
10	j. polski	0,23	0,27	0,18	0,21	0,40	0,17
11	j. polski	0,25	0,29	0,20	0,25	0,45	0,16
12	j. polski	0,28	0,31	0,23	0,28	0,47	0,15
13	j. polski	0,34	0,36	0,27	0,31	0,50	0,13
14	j. polski	0,37	0,41	0,32	0,37	0,61	0,13
15	j. polski	0,42	0,47	0,37	0,41	0,68	0,16
16	j. polski	0,49	0,57	0,45	0,45	0,73	0,16
17	j. polski	0,54	0,67	0,50	0,50	0,74	0,17
18	j. polski	0,54	0,62	0,46	0,44	0,63	0,14
19	j. polski	0,34	0,34	0,26	0,27	0,31	0,10
20	j. polski	0,13	0,09	0,09	0,09	0,07	0,06
21	j. polski	0,02	0,01	0,01	0,01	0,01	0,10
22	j. polski	1,90	2,72	2,43	2,28	3,76	0,90
23	j. polski	1,60	2,20	1,96	1,78	2,80	0,82
24	j. polski	1,46	1,95	1,80	1,56	2,36	0,81
25	j. polski	1,44	1,91	1,80	1,59	2,28	0,85
...

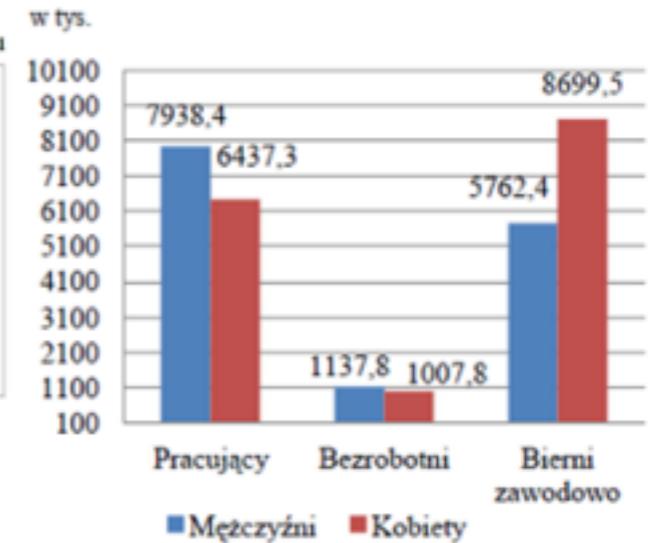
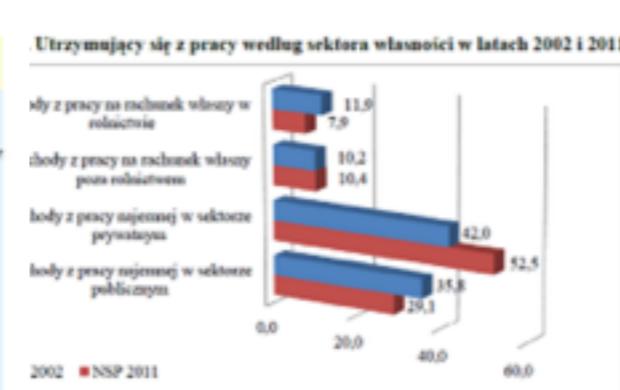
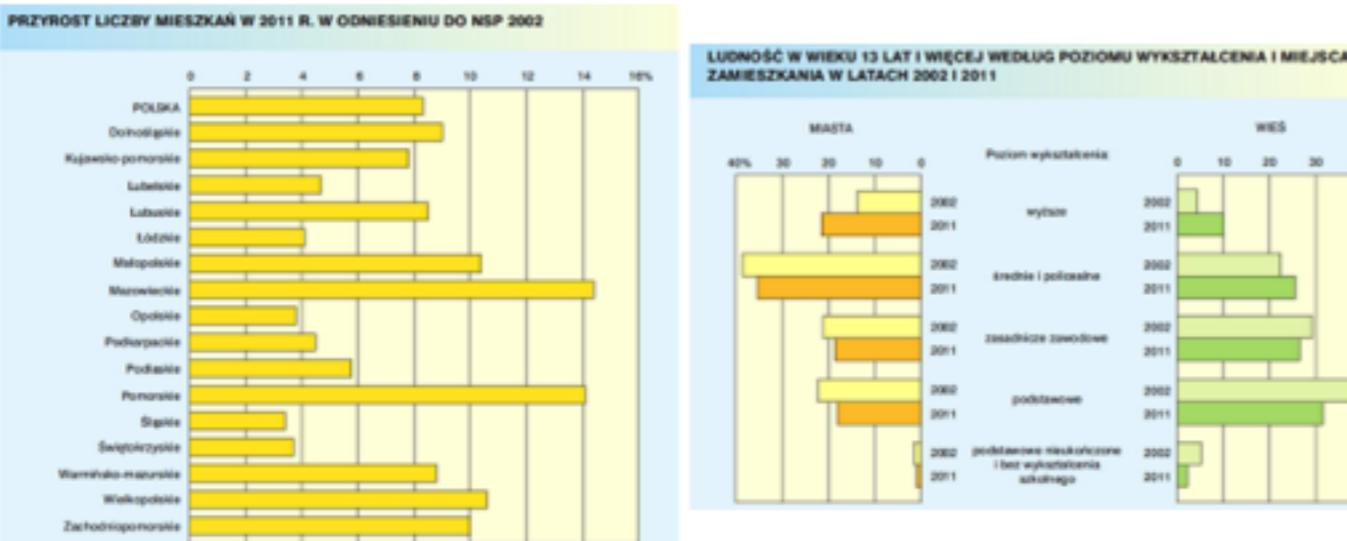
Rozkład liczby punktów na maturze, poziom podstawowy



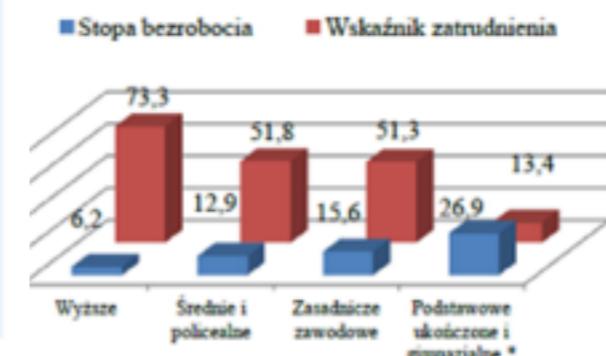
Czy wizualizacja danych jest prosta?

Raporty GUS

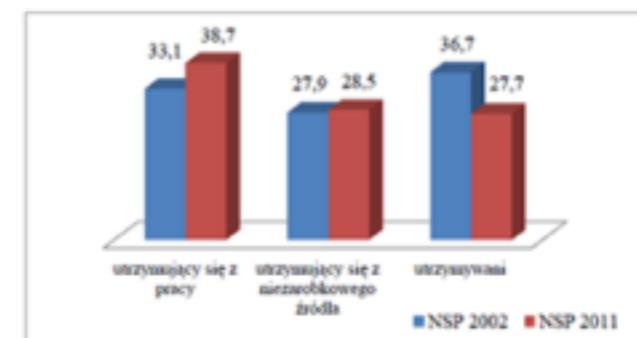
Dlaczego należy uważać z Exceliem?



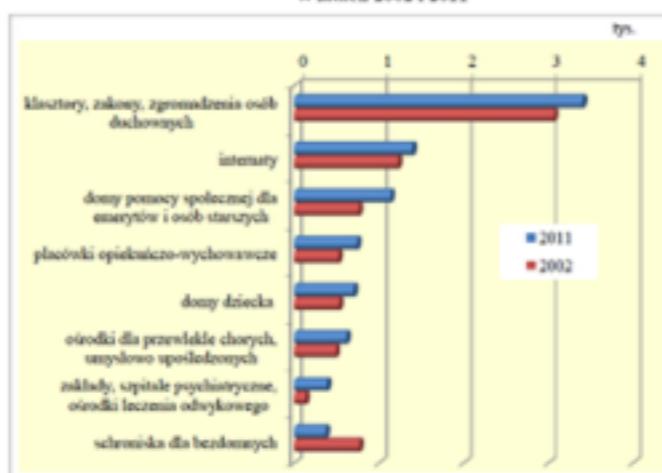
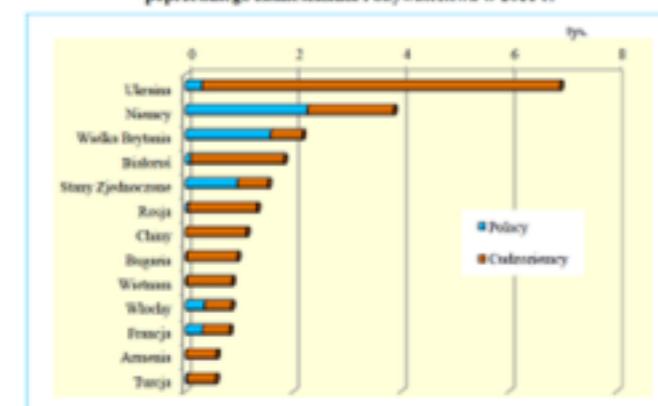
ształcenia w 2011 r.



Wykres 3. Ludność w miastach według głównego źródła utrzymania w latach 2002 i 2011



Wykres 12. Imigranci przebywający w Polsce czasowo powyżej 3 miesięcy według kraju poprzedniego zamieszkania i obywatelstwa w 2011 r.

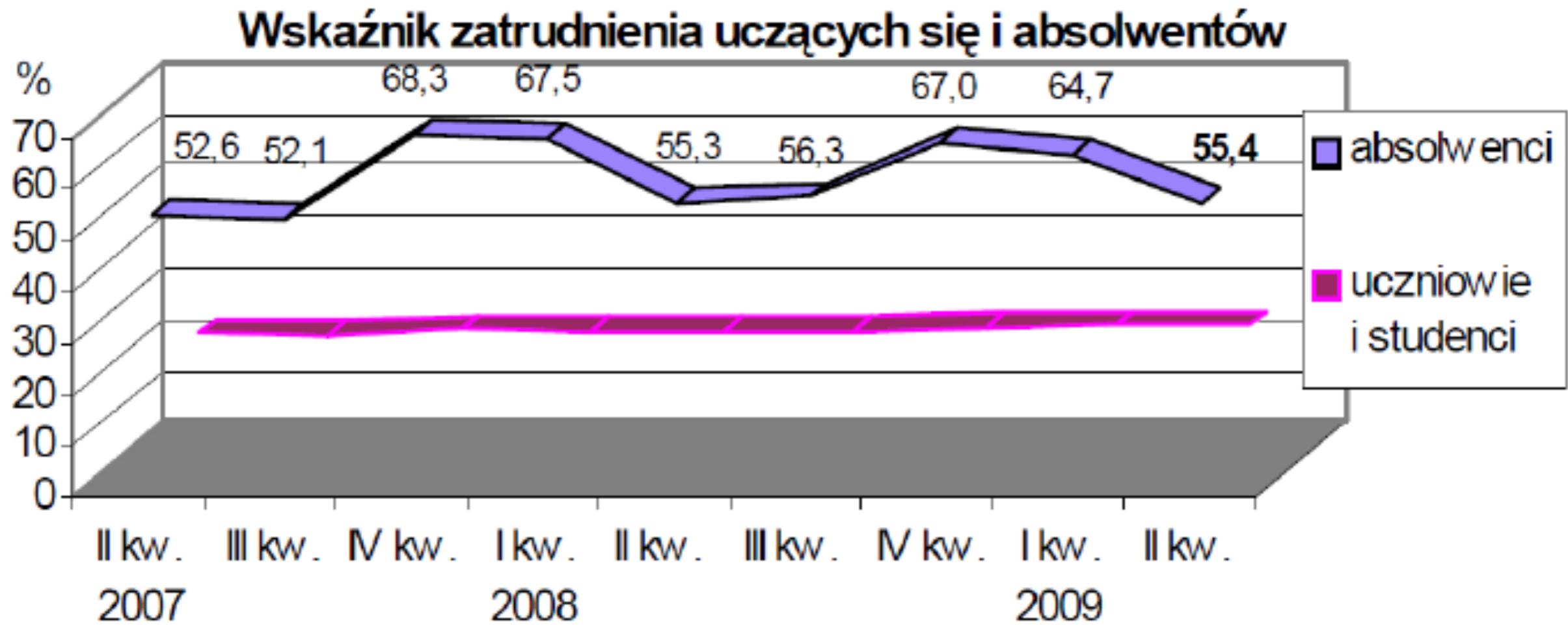


W WIEKU 13 LAT I WIĘCEJ Z WYKSZTAŁCENIEM CO NAJMENЬŚIM W 2002 I 2011



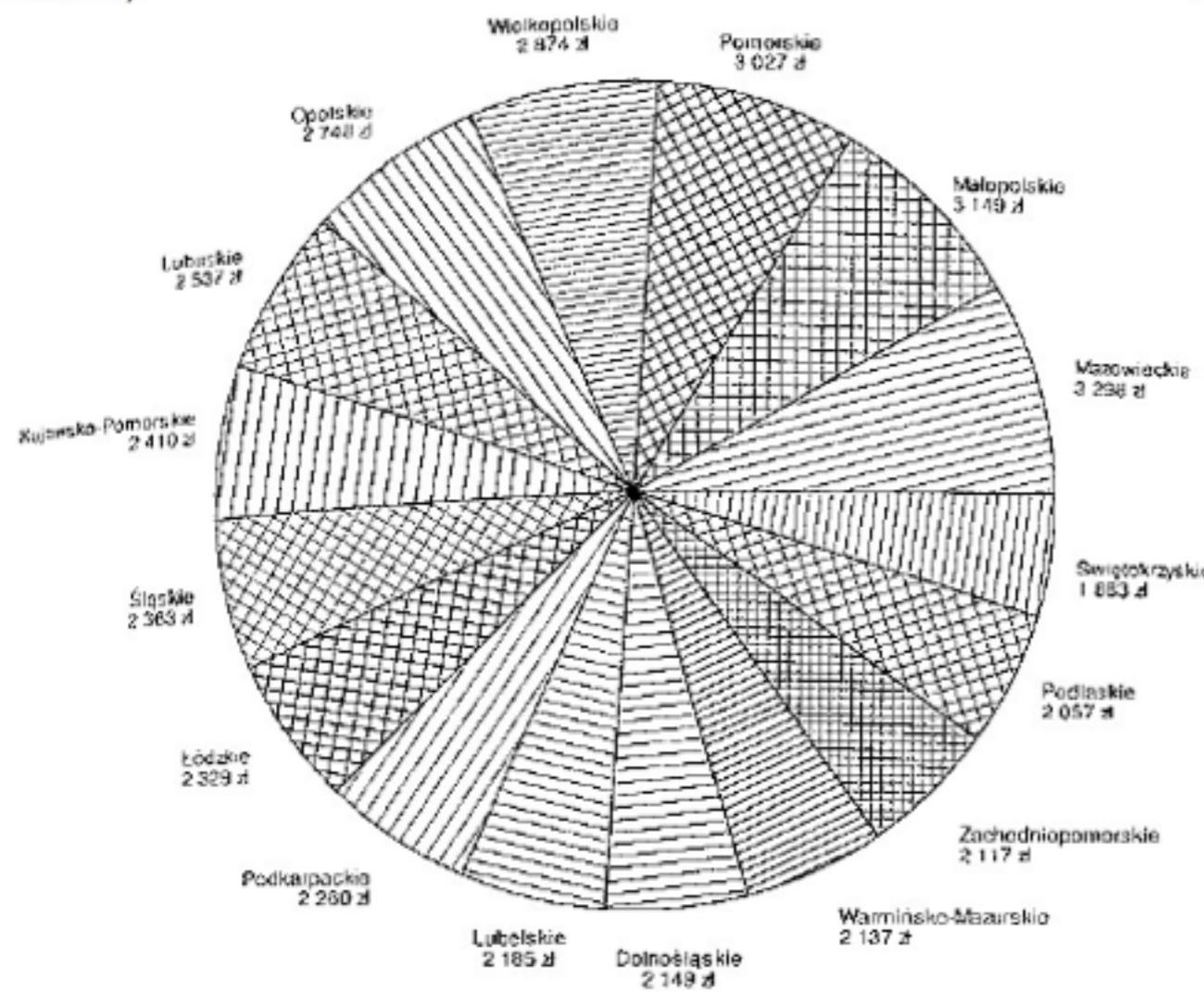
Raporty GUS

Dlaczego należy uważać z Exceliem?

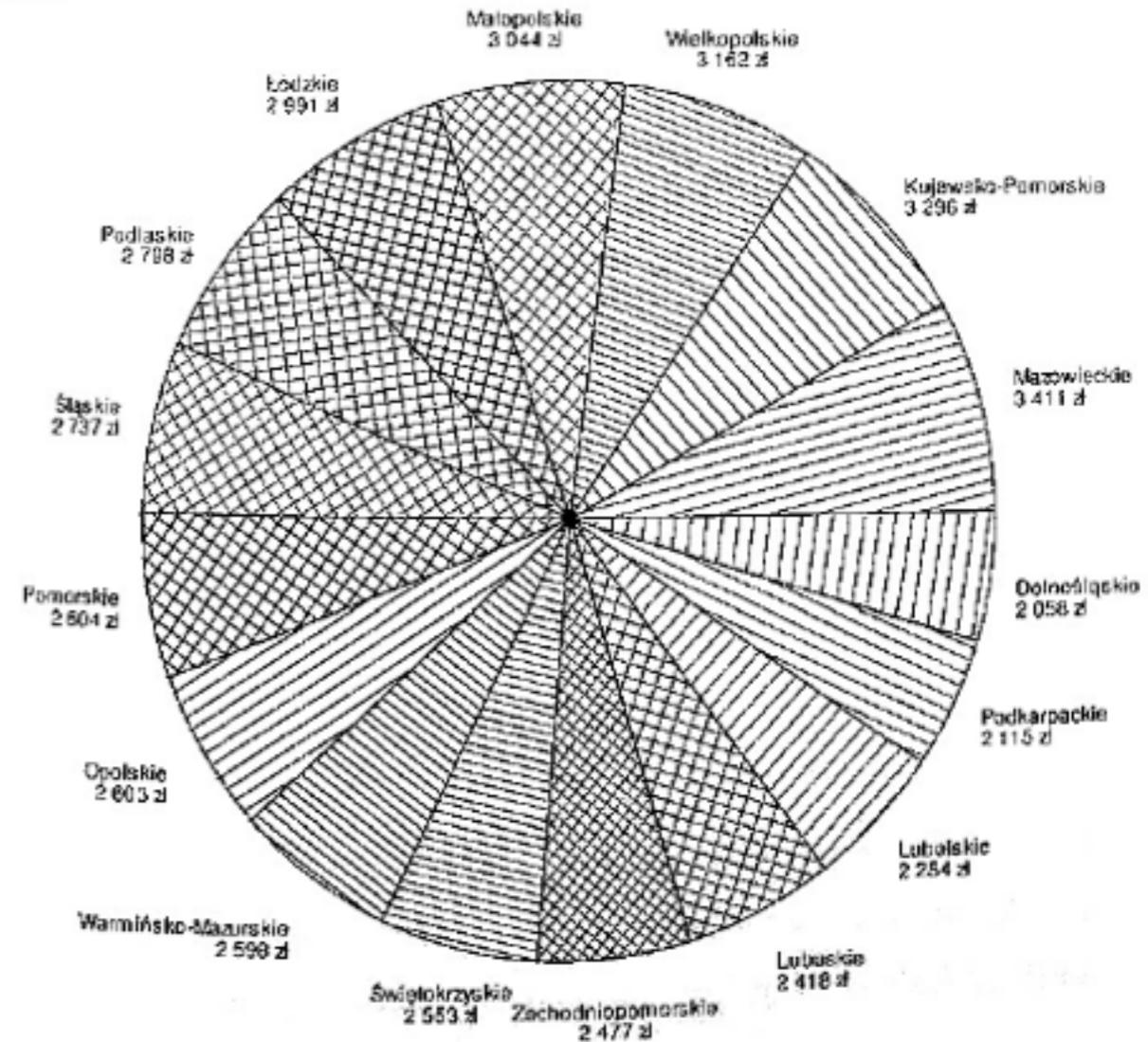


Średni dochód kobiet jest inny niż mężczyzn...

Kobiety



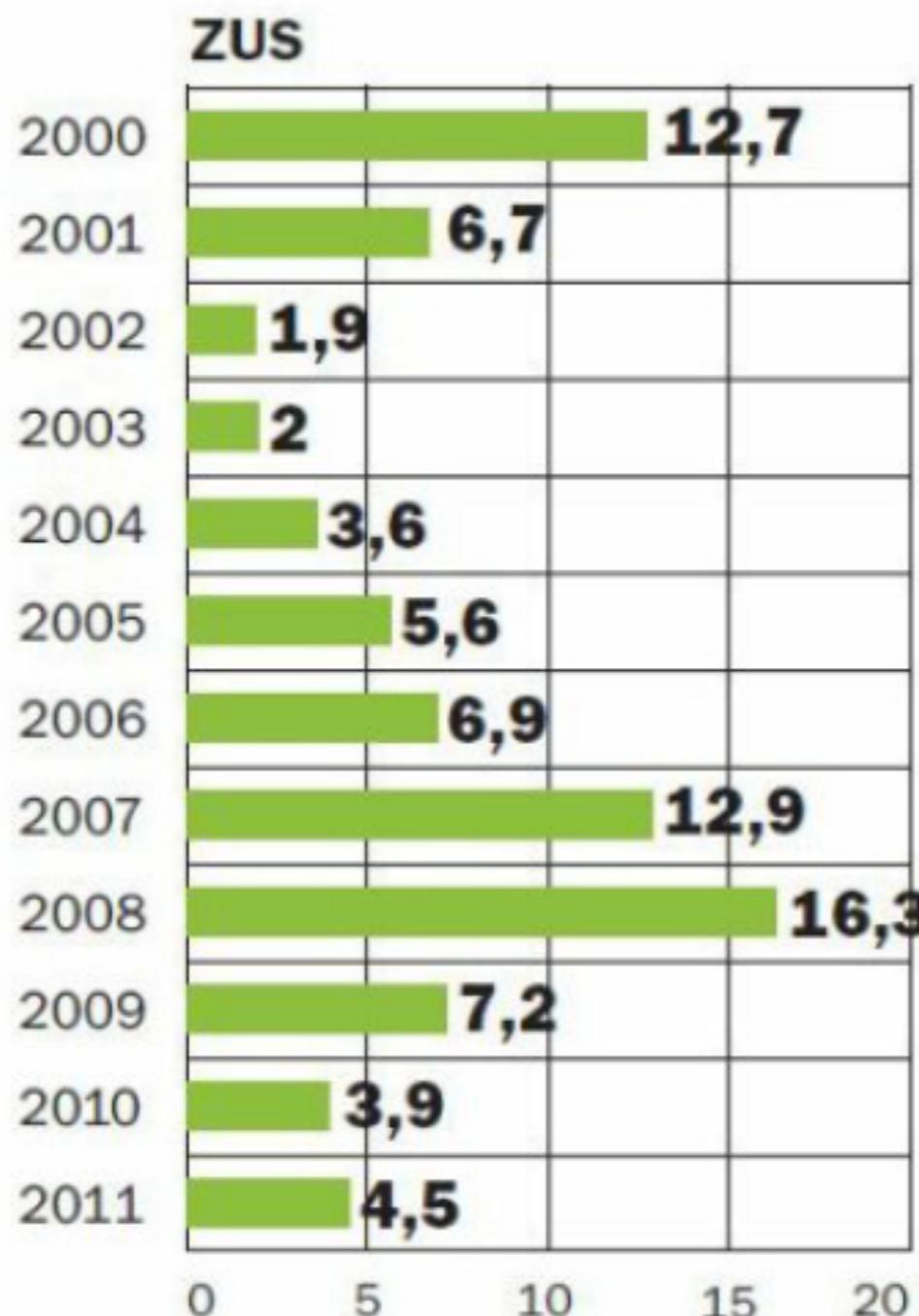
Mężczyźni



Z ZUS można zyskać znacznie więcej! Czyżby?

ILE ZYSKALIŚMY W ZUS I OFE

DANE W PROC.



Znaczny wzrost kosztów użytkowania! Czyżby?

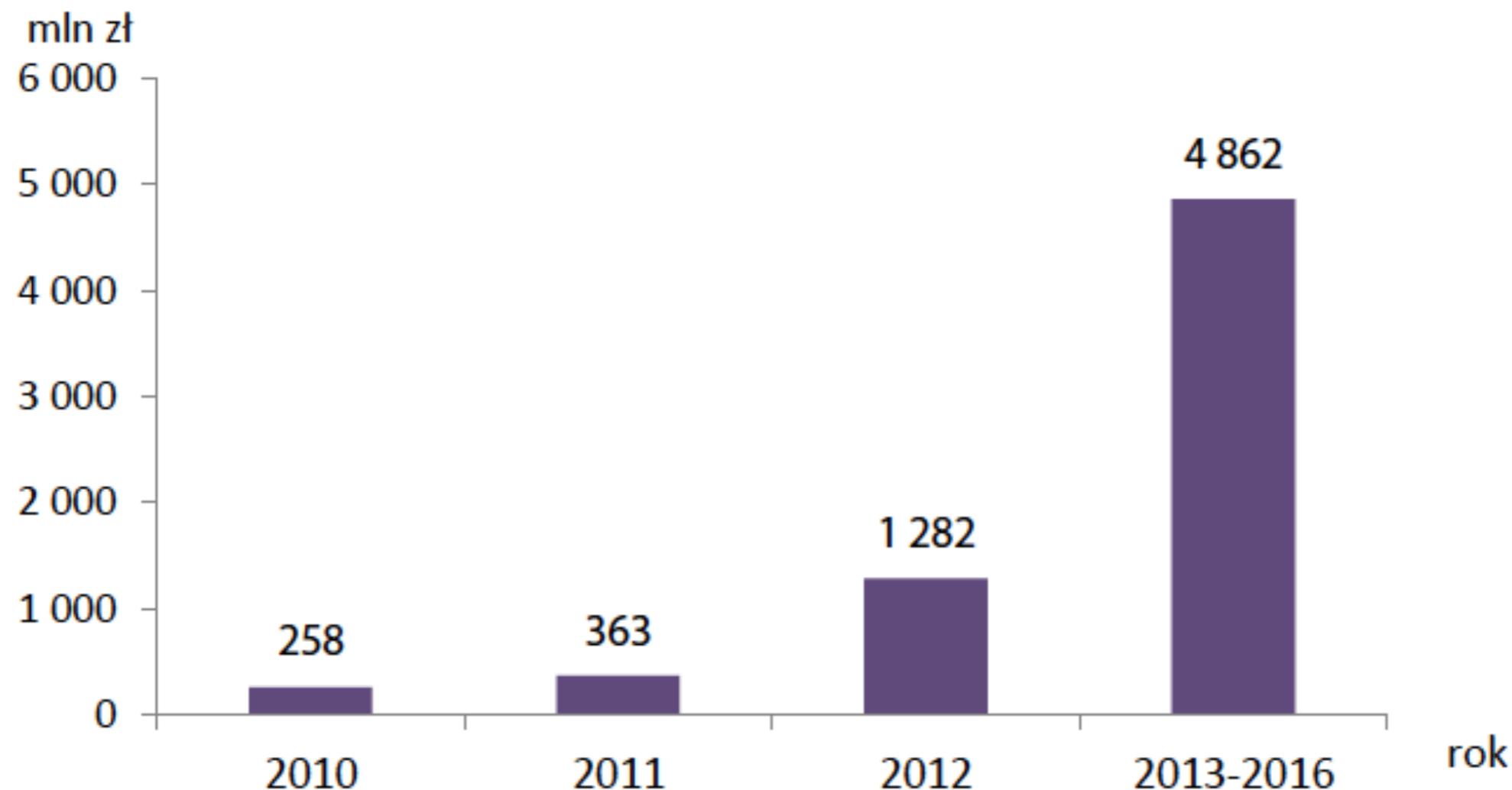
Koszt użytkowania nieruchomości na osobę w gospodarstwie domowym



Źródło: Home Broker, GUS

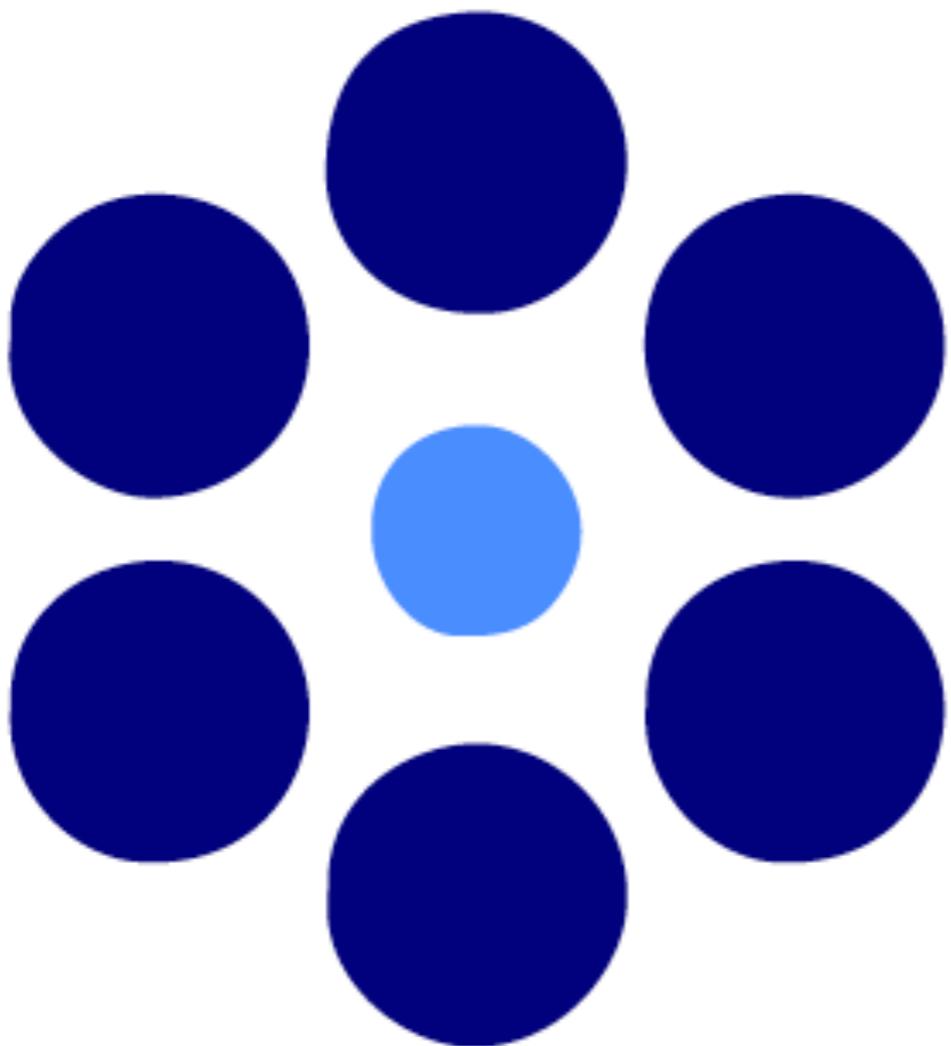
Inwestujemy coraz więcej w B+R! Czyżby?

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

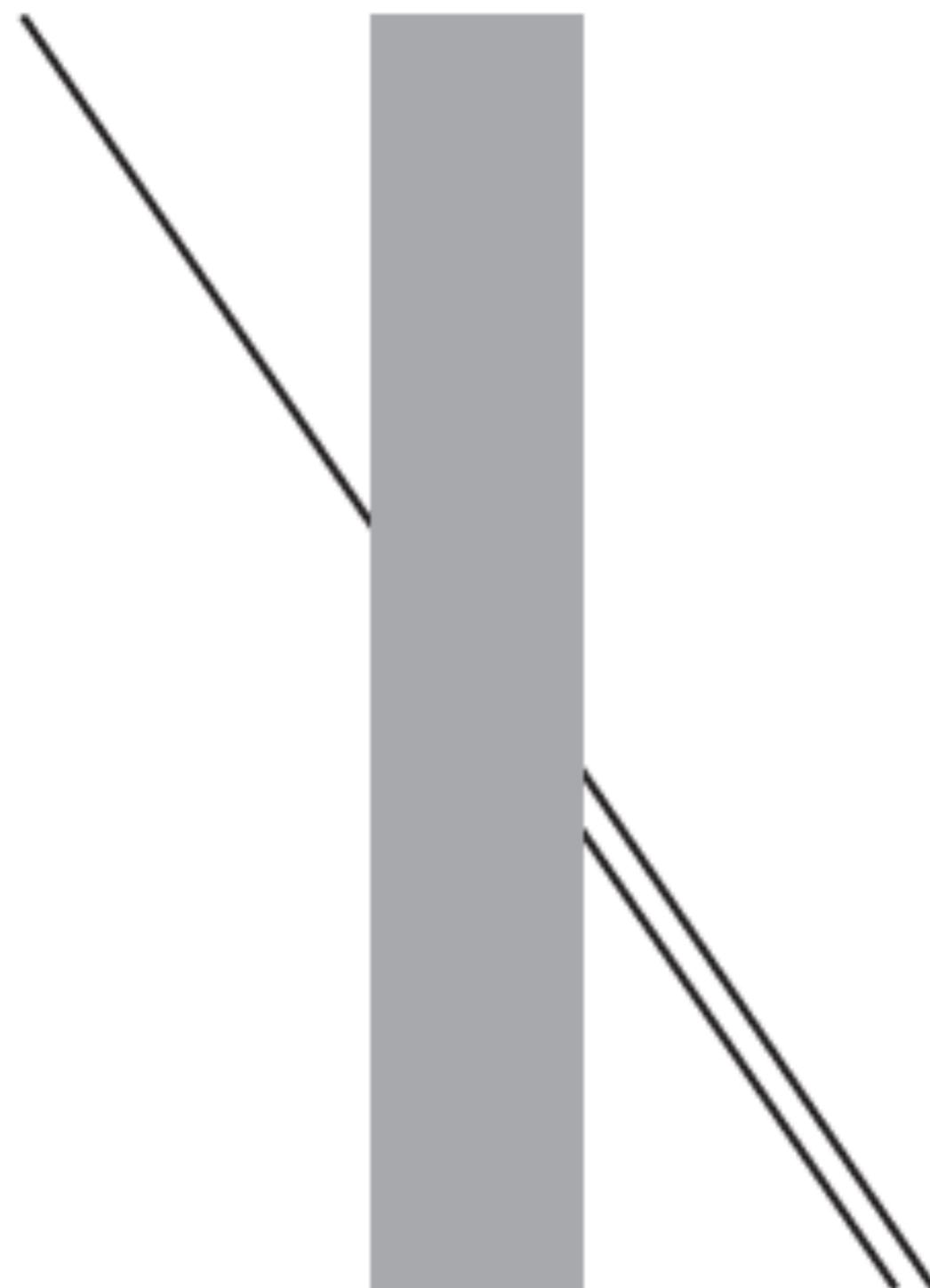


źródło: Narodowe Centrum Badań i Rozwoju

Które z jasnobiebieskich kół jest większe?



Która z prawych linii jest przedłużeniem lewej?

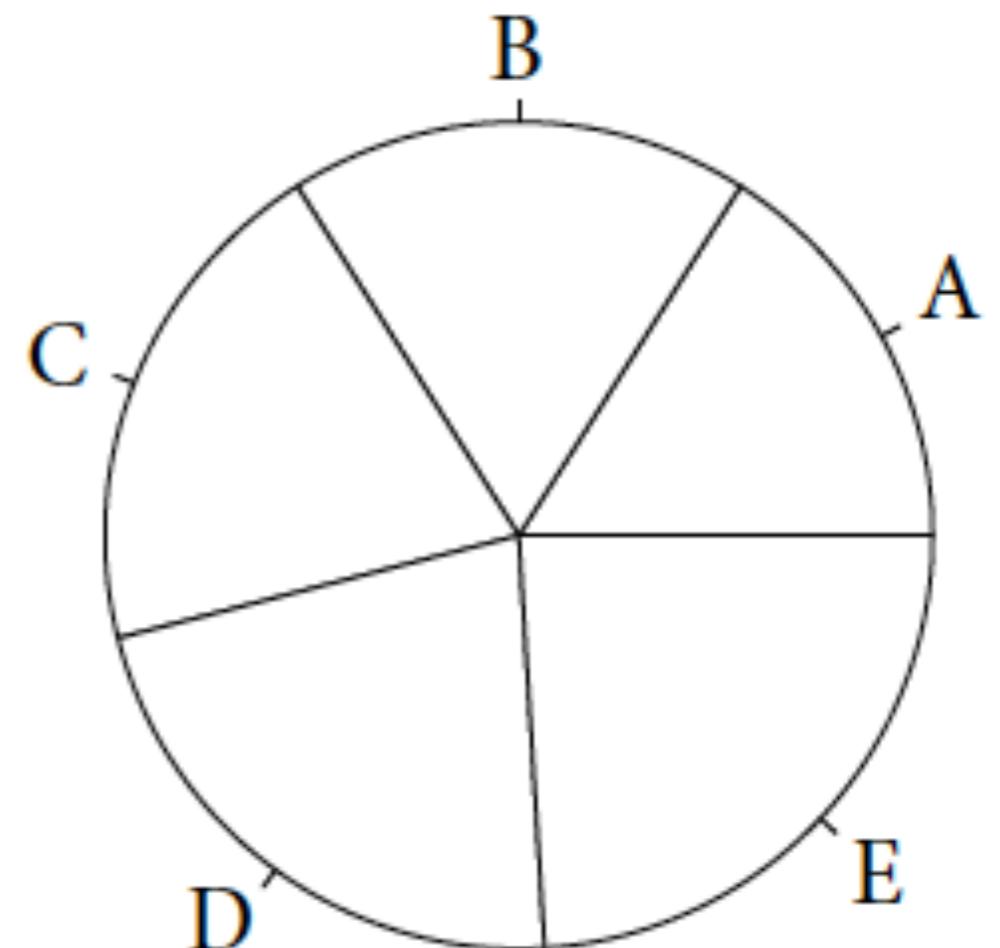
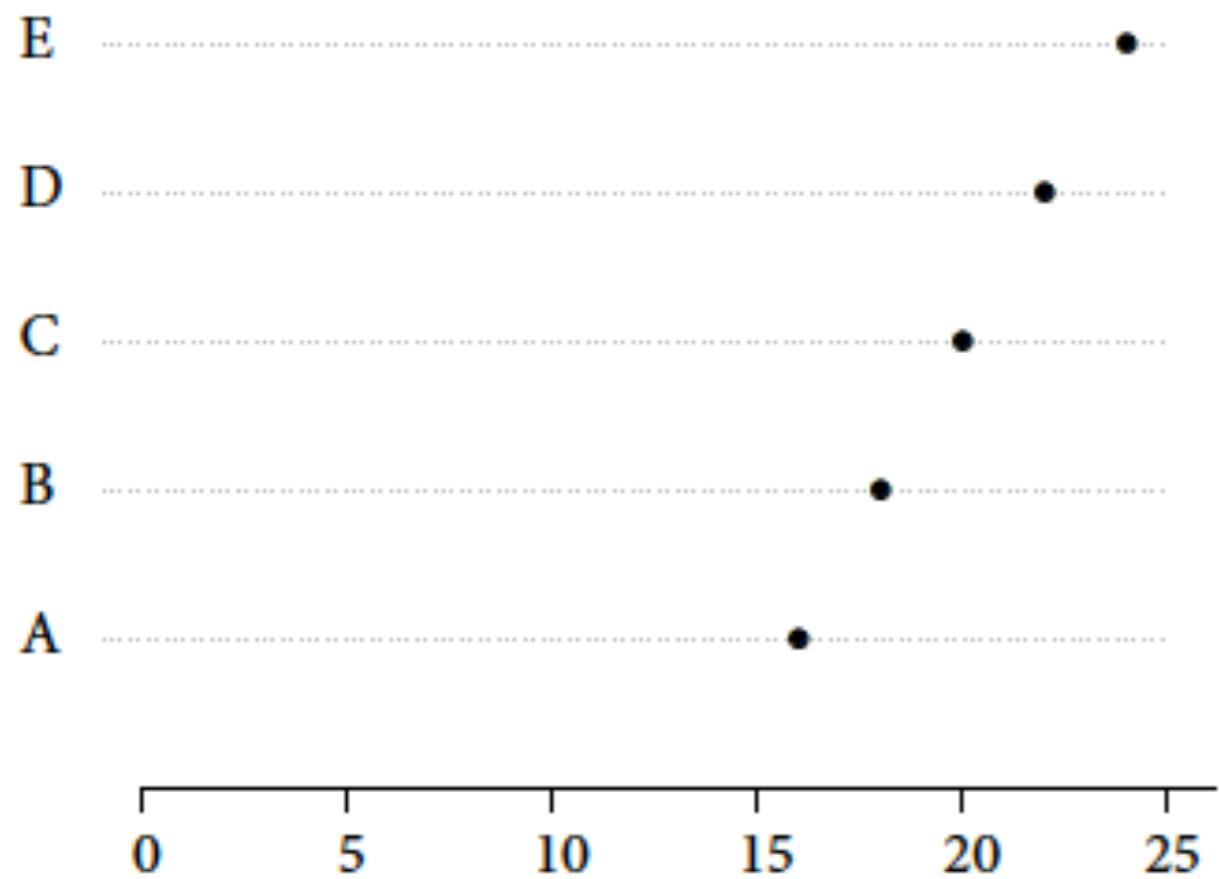




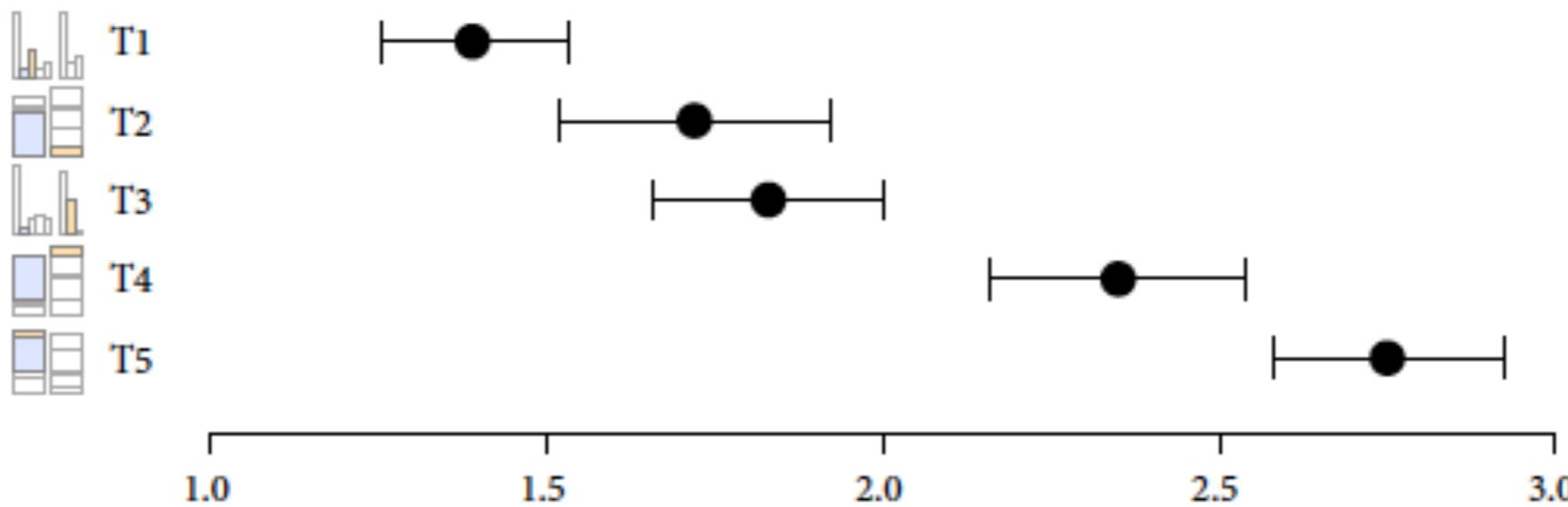
©1999 Daniel J. Simons. All rights reserved.
Image may not be distributed or posted online without written permission

A więc jak żyć?

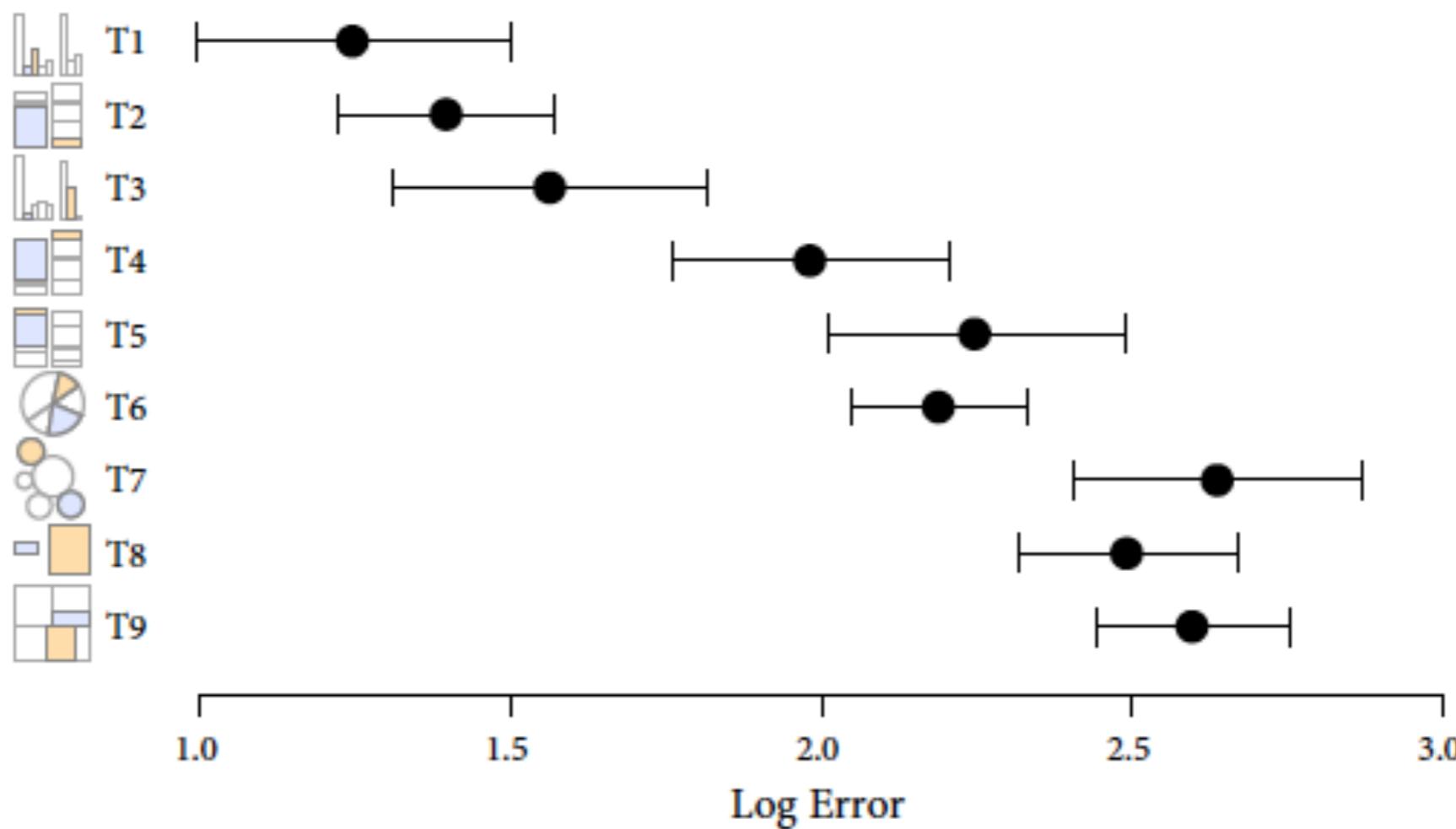
Zagadka 3: Pola, kąty, pozycje czy długości?



Cleveland & McGill's Results



Crowdsourced Results

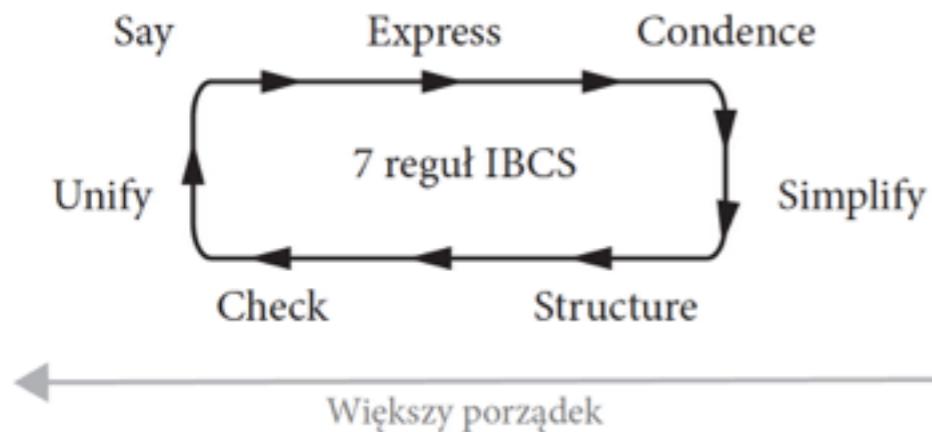


The International Business Communication Standards

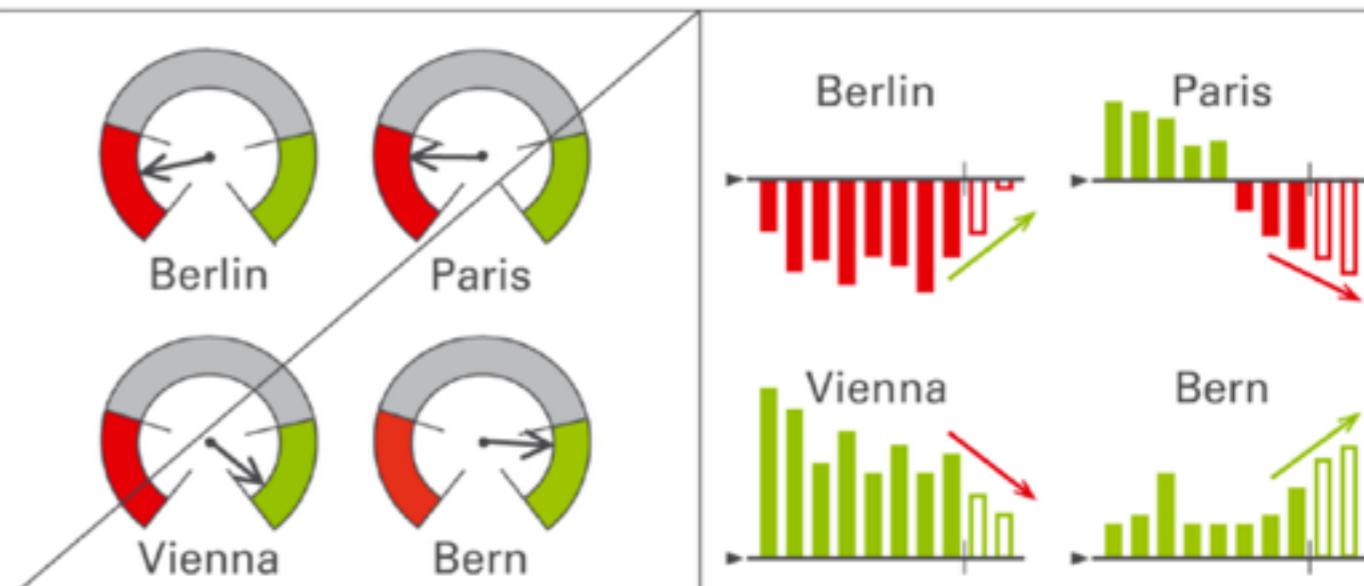
<http://www.ibcs-a.org/>

Say	Przekazuj konkretny, istotny komunikat.
Unify	Stosuj jednolite oznaczenia i wyróżnienia.
Condence	Zwiększać ilość informacji, dbaj o treściwość.
Check	Zapewnij wizualność spójność wykresów i tabel.
Express	Wybierz odpowiednią formę (opis, tabelę, wykres) dla prezentowanych danych.
Simplify	Unikaj bałaganu i zbędnych ozdobników.
Structure	Zadbaj o czytelną strukturę.

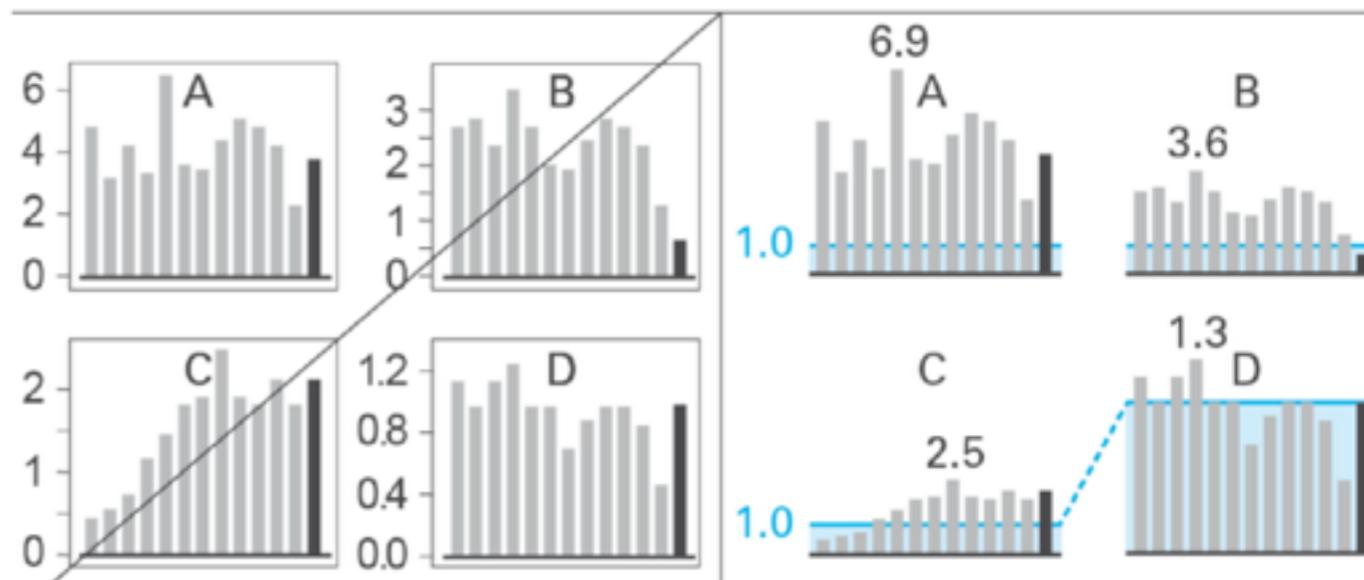
Więcej treści



EX 2.2 Replace gauges, speedometers



UN 5.2 Unify scaling indicators



iHealth, Punkty, łamana czy słupki?

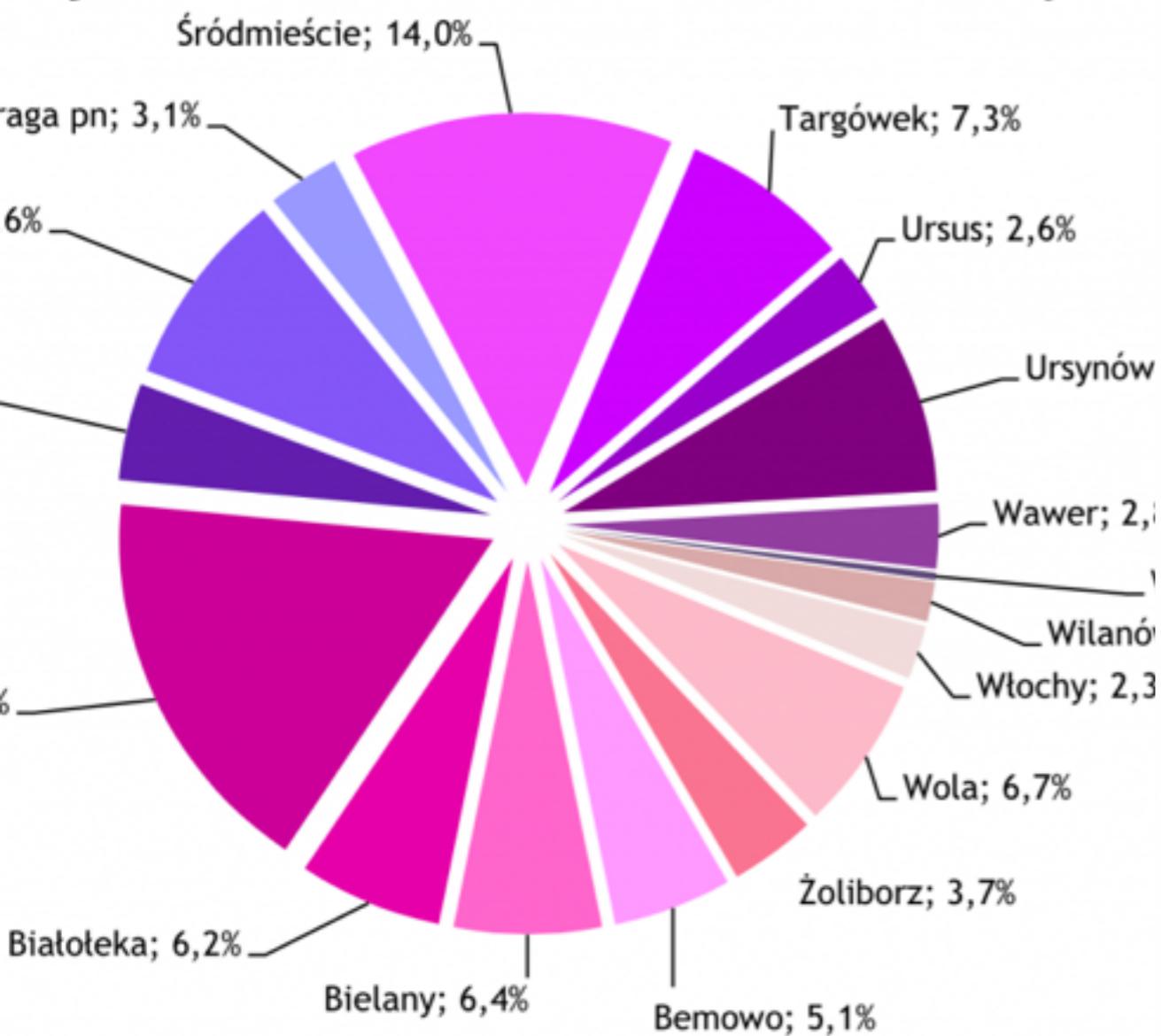


Ok, w takim razie używamy słupków.

A jak dobierać kolory?

Nie ma róży bez ognia

Miejsce zamieszkania - dzielnice Warszawy



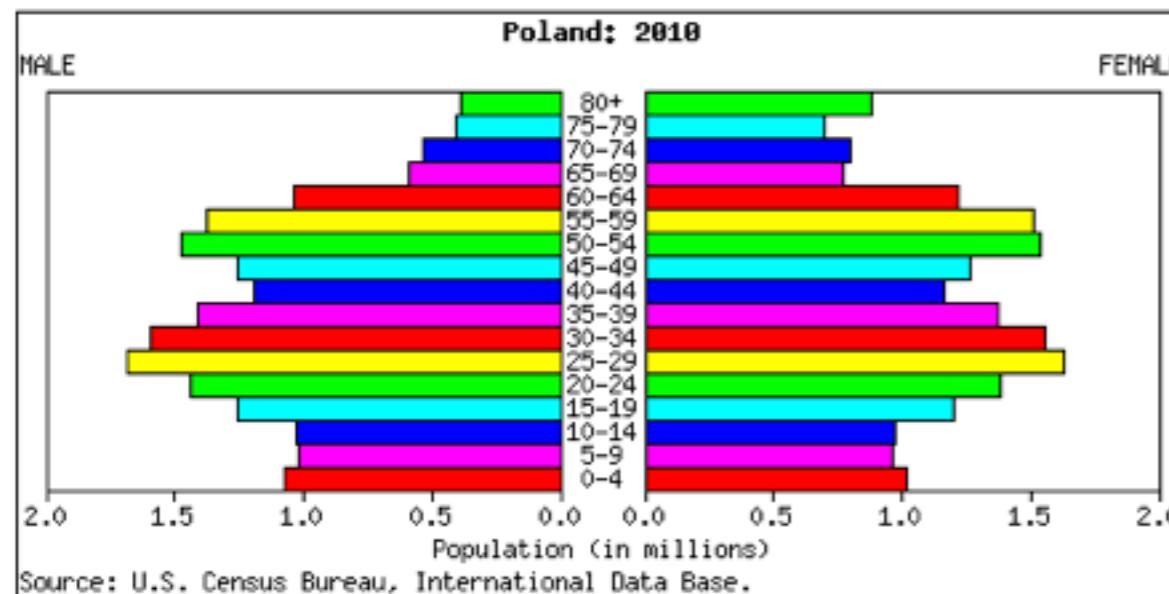
ŚRODOWISKOWE DOMY SAMOPOMOCY
z tablicy nr 7-2 kolumna nr 10; wiersz POLSKA: 728.



Nie ma róžu bez ognia

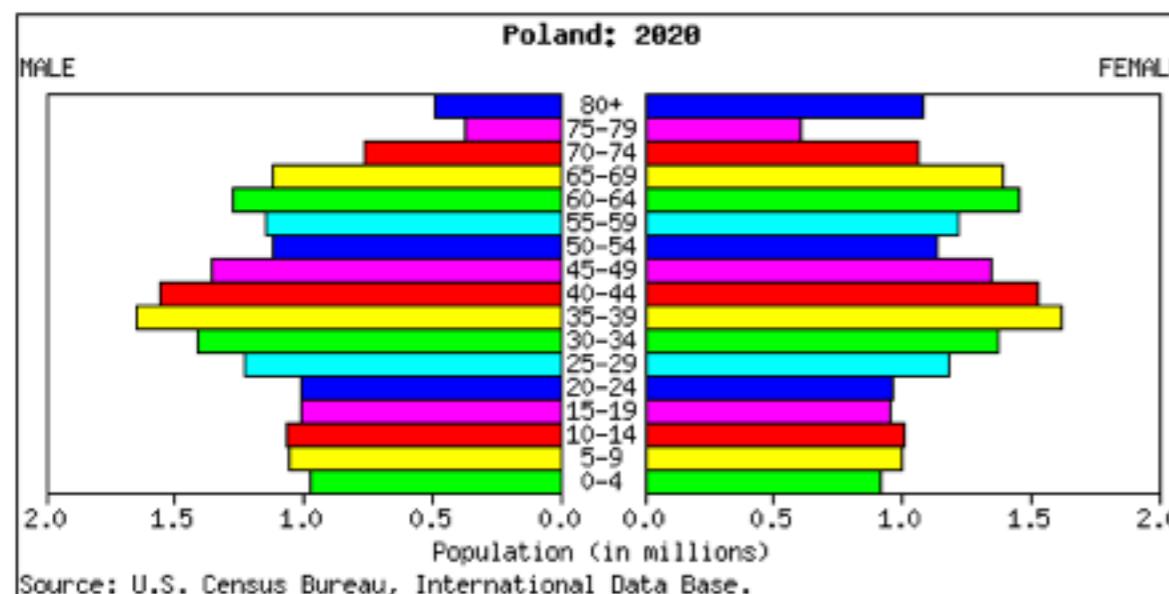
Poland Population Pyramid for 2010

Age and sex distribution for the year 2010:



Poland Population Pyramid for 2020

Predicted age and sex distribution for the year 2020:



OPINION

Obama's Divided Nation

With 66 of states Obama presides over, America more divided than any time in 50 years that was riven by racial lines gathering in 2008 to elect its president. That president has four years dividing the basis of economy. The campaign revealed no evidence that Mr. Obama will close the chasm he has created between his voters and those he attacked and vilified.

It may be true that Mitt Romney failed to respond



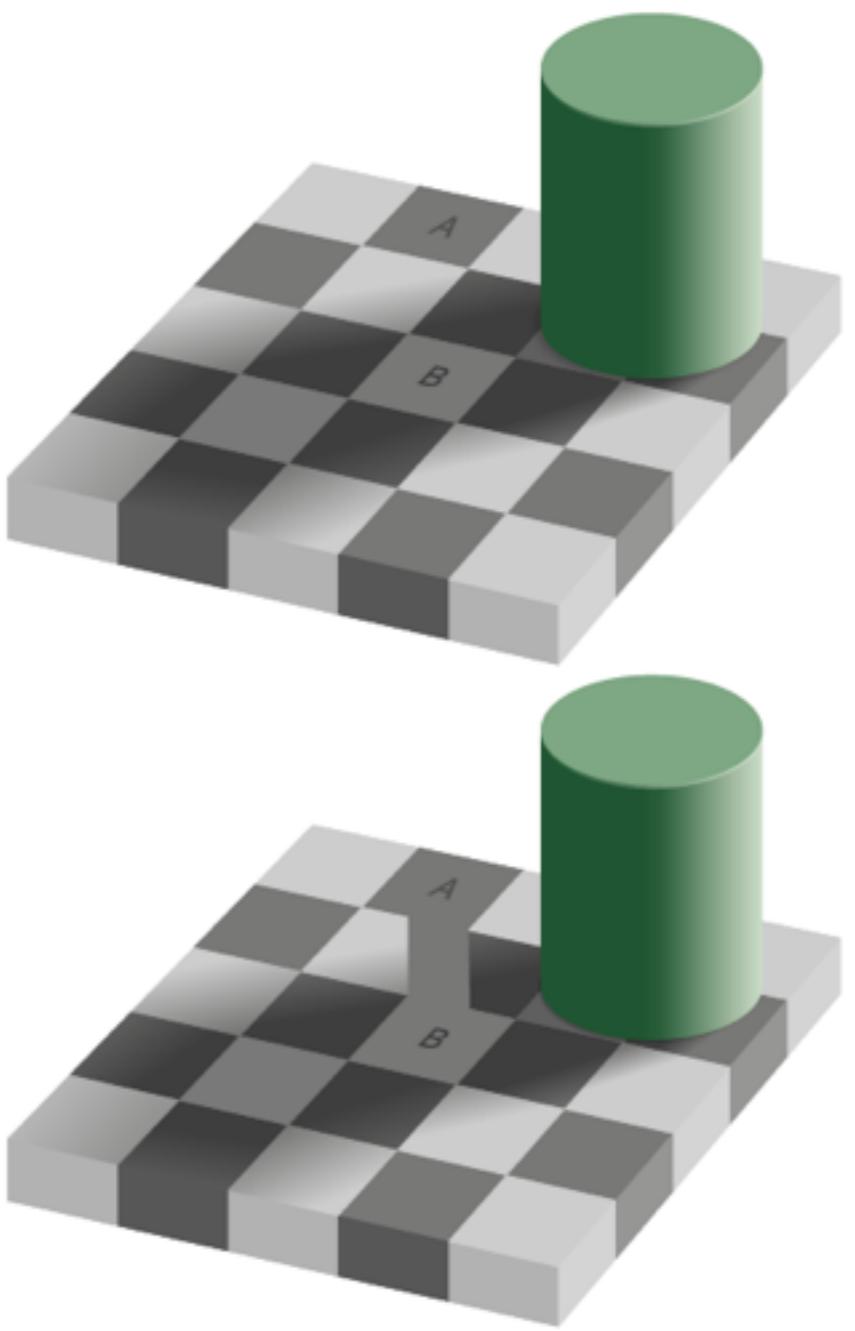
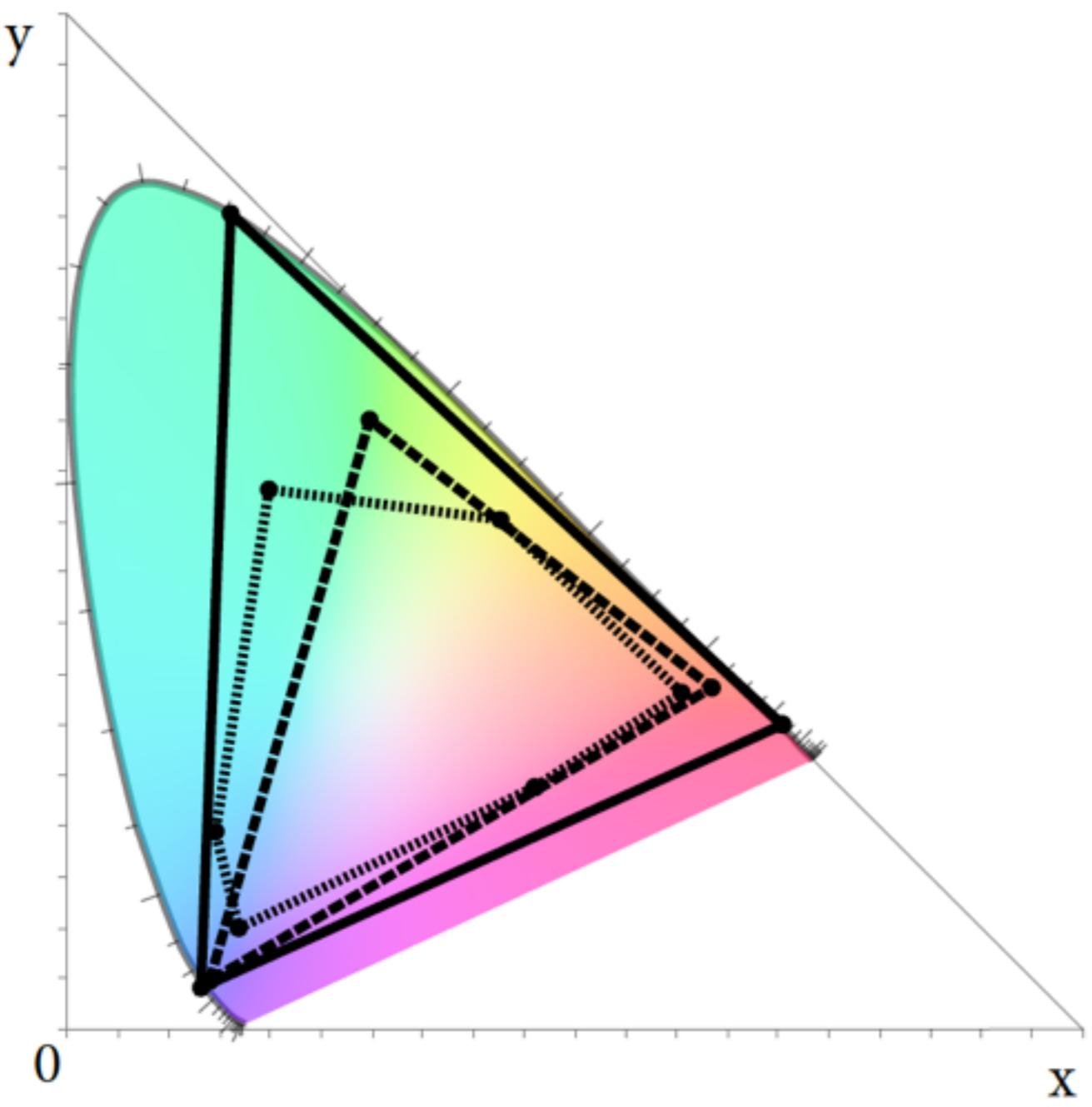
Source: AP

drawn attention to what hap-

Obama spokesman replied: "It's right and you're wrong. It

problem with pols, verbally facile as Mitt Romney is, that in crunch time, reverts to No. 1. Evidence that 9% of the electorate who to vote for just Tuesday; and among them, 42% said Mr. Obama's Sandy response—tie photo-op—was a factor. Of those, 50% voted for Mr. Christie is a politico who is

Yes, Republicans across two presidents that there are some who are how crudely issue like immigration. Blowing up the if you thought day's results

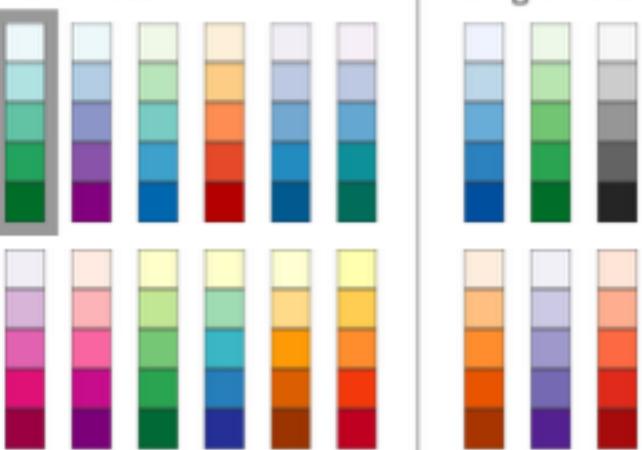


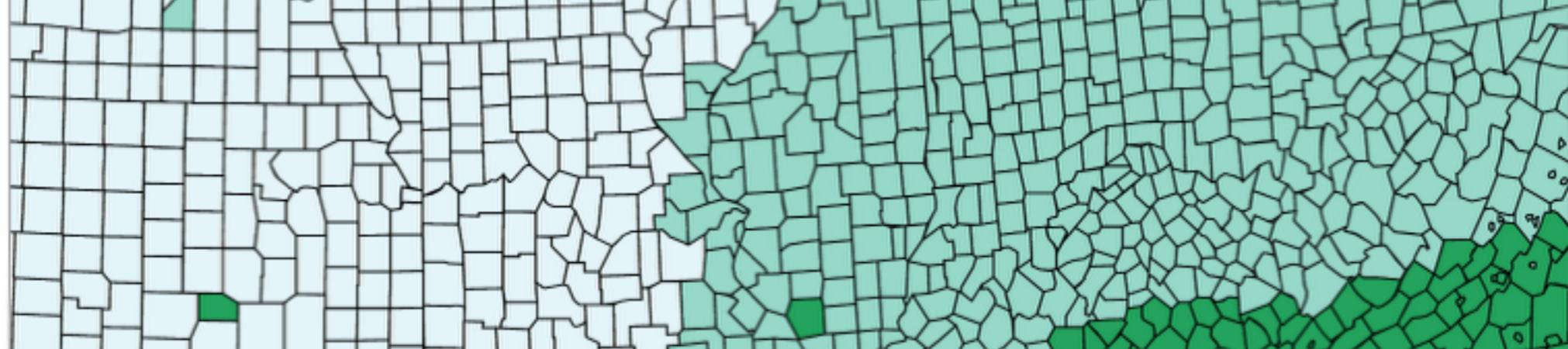
Kim jest Cynthia Brewer?

Number of data classes: 3 how to use | updates | downloads | credits

Nature of your data: sequential diverging qualitative

Pick a color scheme:

Multi-hue: 

Single hue: 

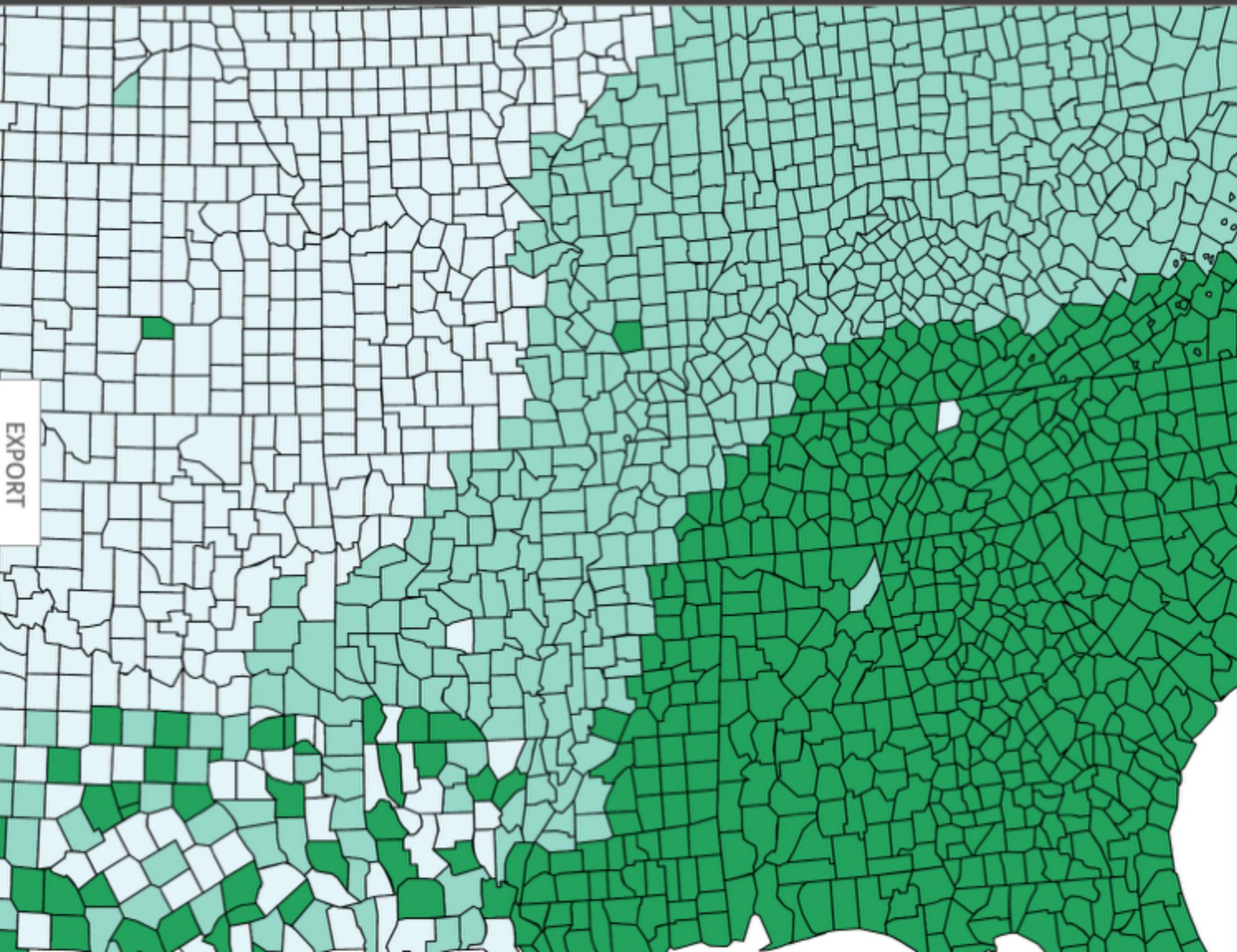
Only show: 3-class BuGn

colorblind safe print friendly photocopy safe

Context: roads cities borders 

Background: solid color terrain

#e5f5f9
#99d8c9
#2ca25f



COLORBREWER 2.0
color advice for cartography

Colours In Culture



A Western / American	1 Anger
B Japanese	2 Art / Creativity
C Hindu	3 Authority
D Native American	4 Bad Luck
E Chinese	5 Balance
F Asian	6 Beauty
G Eastern European	7 Calm
H Muslim	8 Celebration
I African	9 Children
J South American	10 Cold

19 Desire	37 Growth
20 Earth	38 Happiness
21 Energy	39 Healing
22 Eros	40 Healthy
23 Etc	41 Heat
24 Evil	42 Heaven
25 Excellence	43 Holiness
26 Far	44 Illness
27 Fer	45 Insight
28 Fer	46 Intelligence
29 Flare	47 Intuition
30 Fre	48 Religion
31 Fri	49 Jealousy
32 Fur	50 Joy
33 Go	51 Learning
34 Go	52 Life
35 Go	53 Love
36 Gra	54 Loyalty
37 Roy	55 Luxury
38 Sel	56 Marriage
39 Stre	57 Modesty
40 Sty	58 Money
41 Sud	59 Mourning
42 Tro	60 Mystery
43 Tru	61 Nature
44 Un	62 Passion
45 Peac	63 Peace
46 Pen	64 Penance
47 Pow	65 Power
48 Perso	66 Personal power
49 Pur	67 Purity
50 Rad	68 Radicalism
51 Rat	69 Rational
52 Lif	70 Reliable
53 Lov	71 Repels Evil
54 Loy	72 Respect

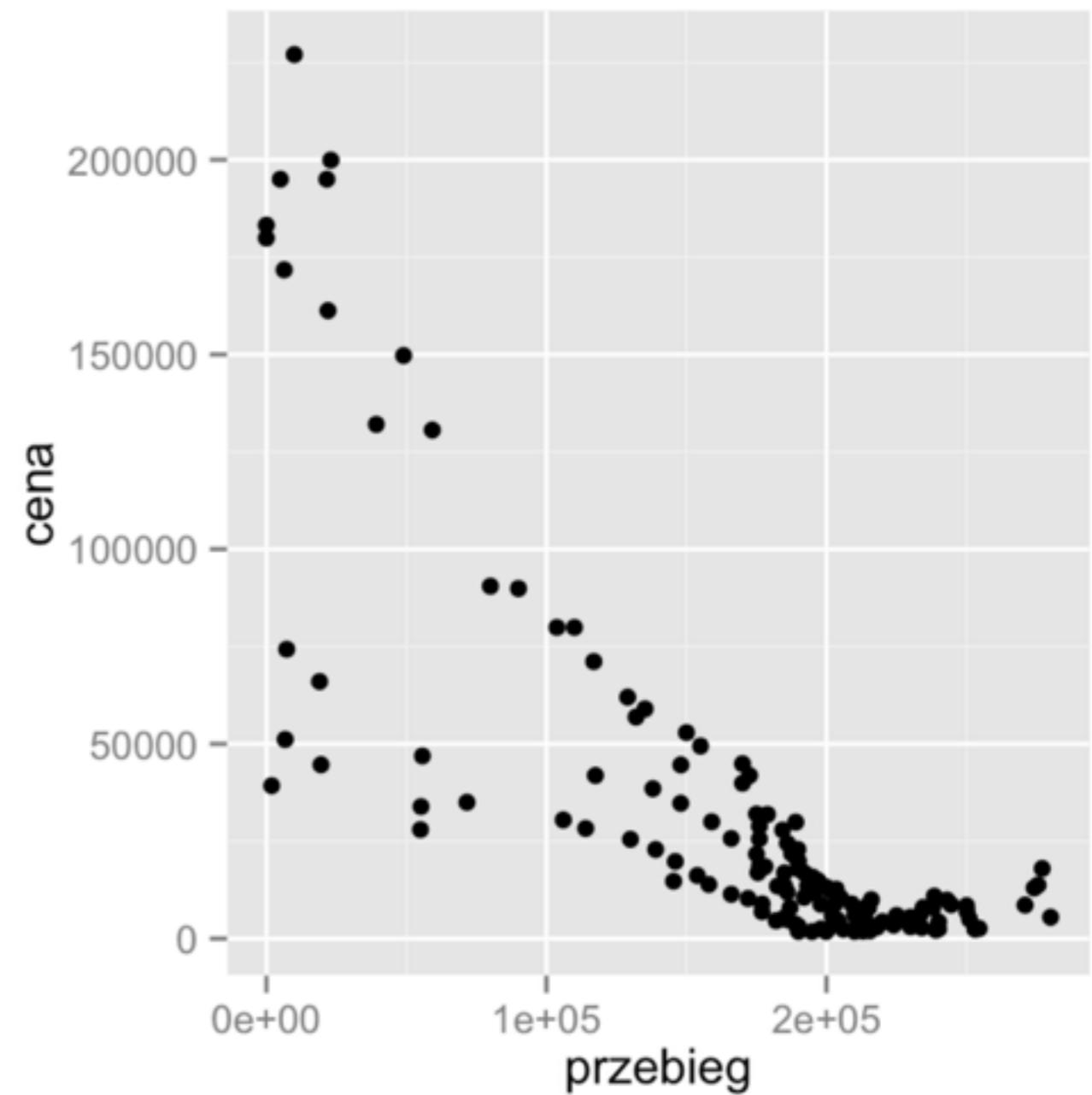
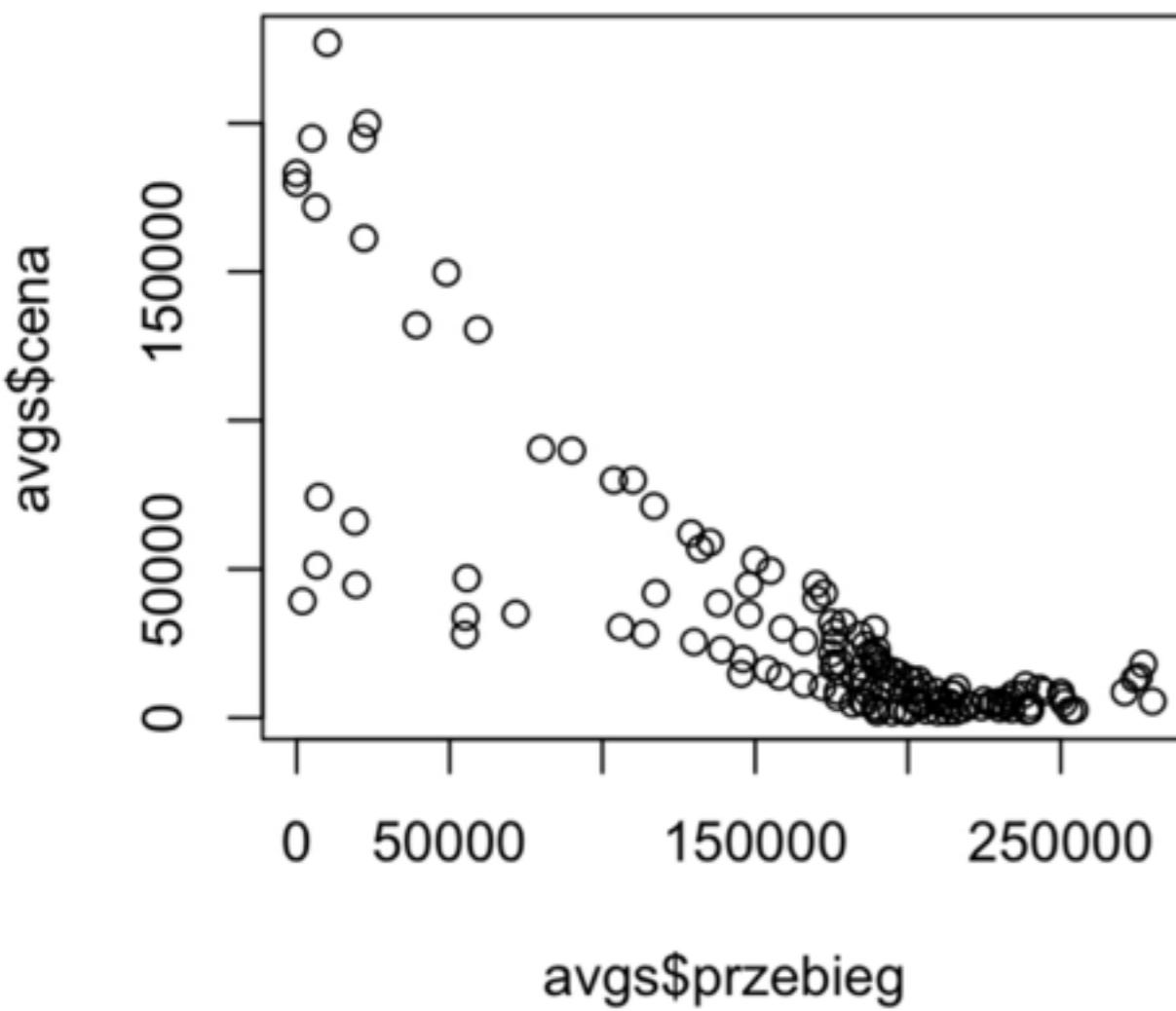
█ Yellow █ Grey
█ Gold █ Silver

Czas na R



Trzy podejścia do tworzenia statycznej grafiki

```
plot(avgs$przebieg, avgs$cena)          # base graphics  
ggplot(avgs, aes(x=przebieg, y=cena)) + # ggplot graphics  
  geom_point()
```

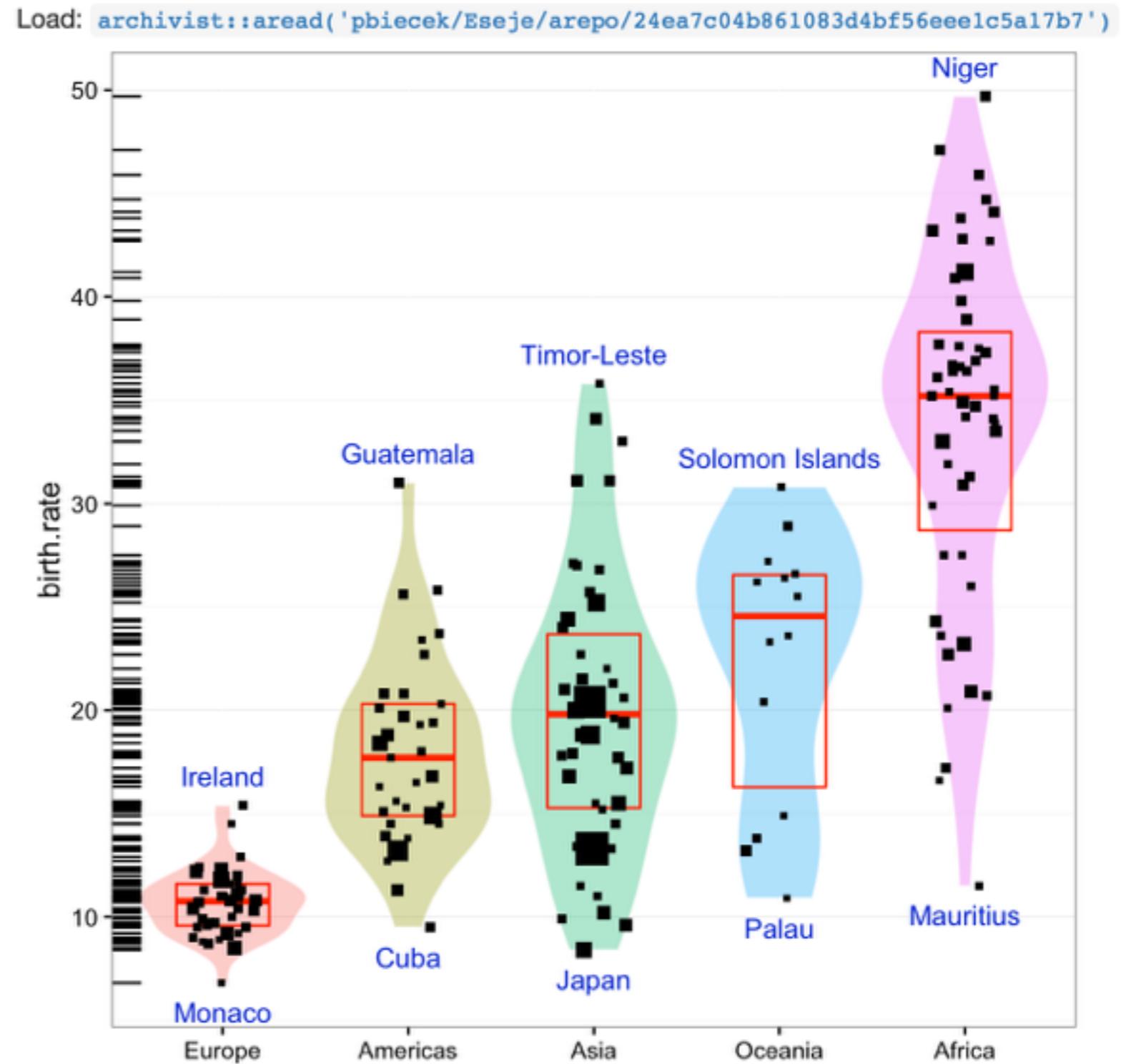


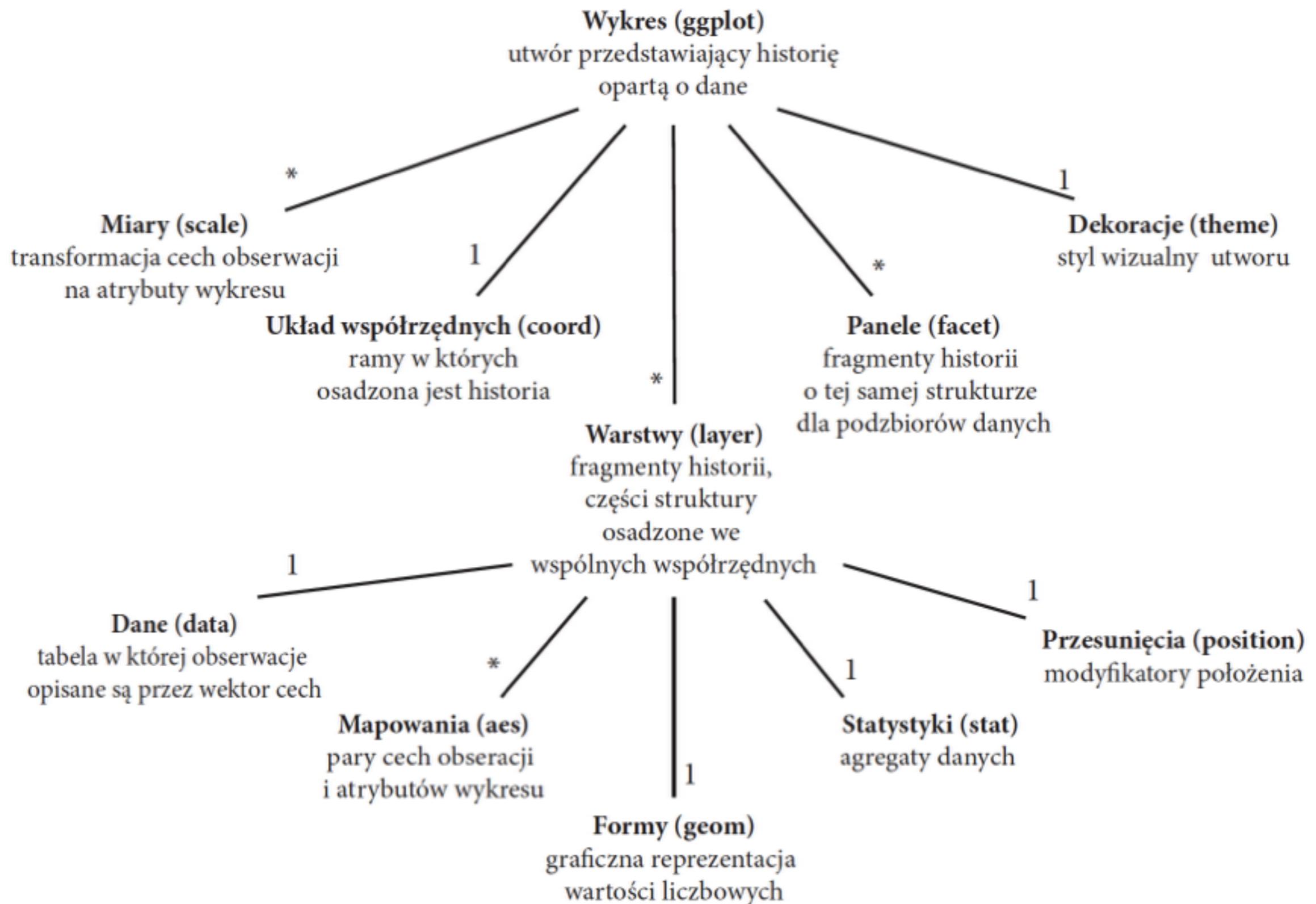
A więc dlaczego ggplot2?

Pakiety:
ggplot2
archivist
ggThemeAssistant

Showcase:

[https://rawgit.com/
pbiecek/Eseje/master/
GamatykaGrafiki.html](https://rawgit.com/pbiecek/Eseje/master/GamatykaGrafiki.html)



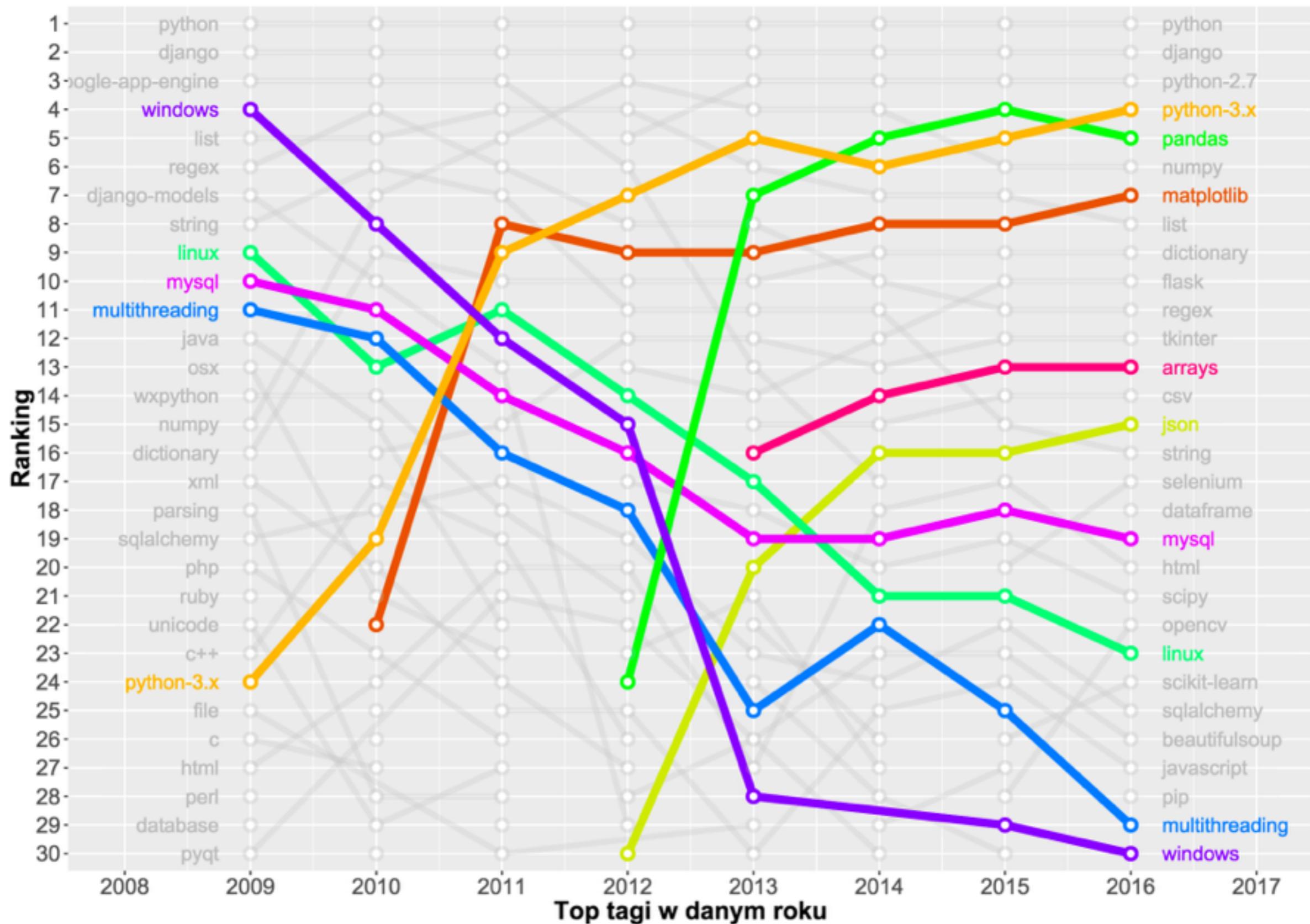


datahero.tech

Analiza rynku pracy dla osób pracujących z danymi:

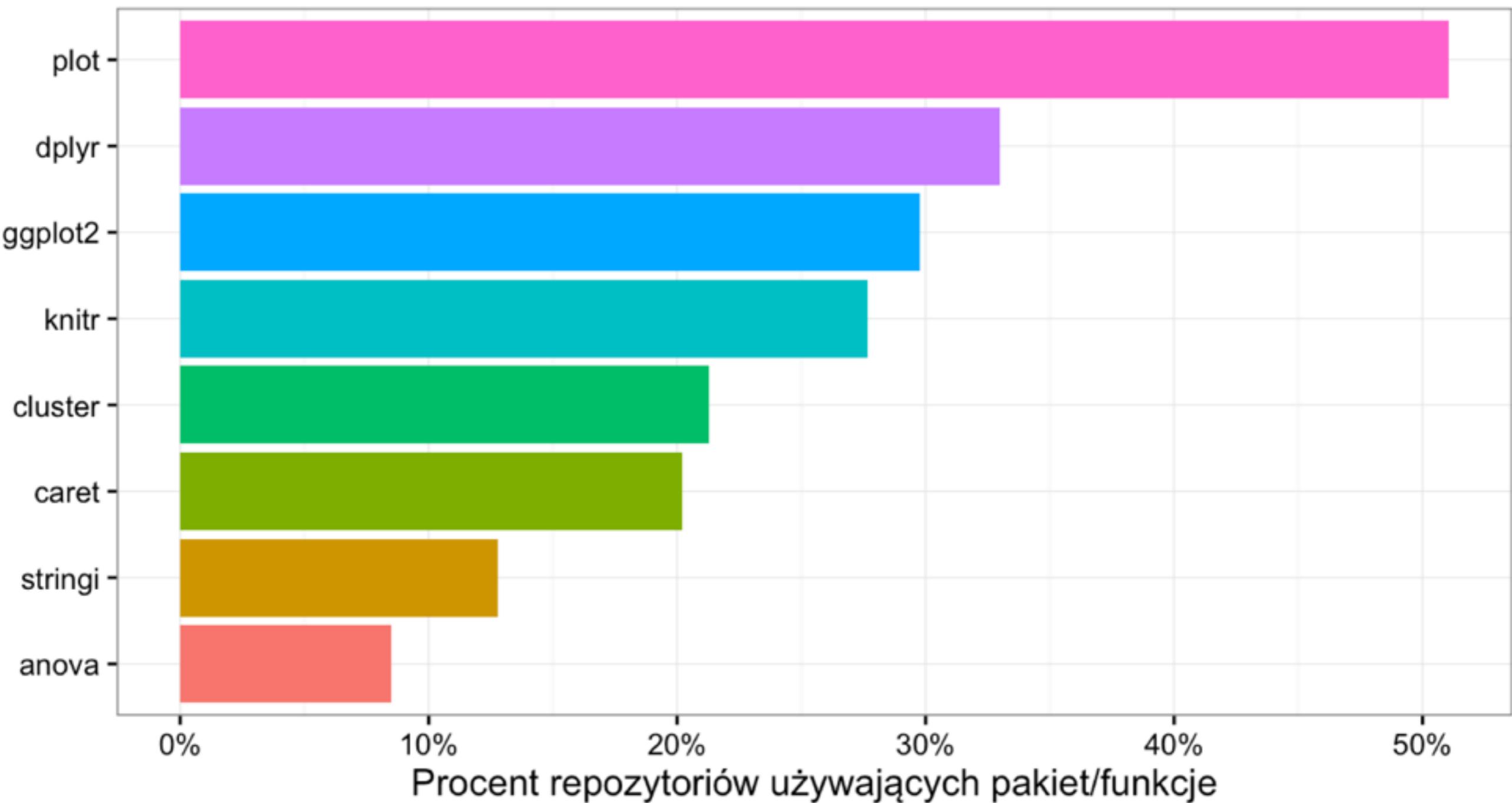
- Olga Mierzwa @ DataHero - analiza StackOverflow
- Konrad Więcko, Krzysztof Słomczyński @ MI2 - zbieranie i analiza pracuj.pl

StackOverflow:



GitHub:

Skanujemy repozytoria R i sprawdzamy kto i jak często korzysta z określonych funkcji/funkcjonalności.



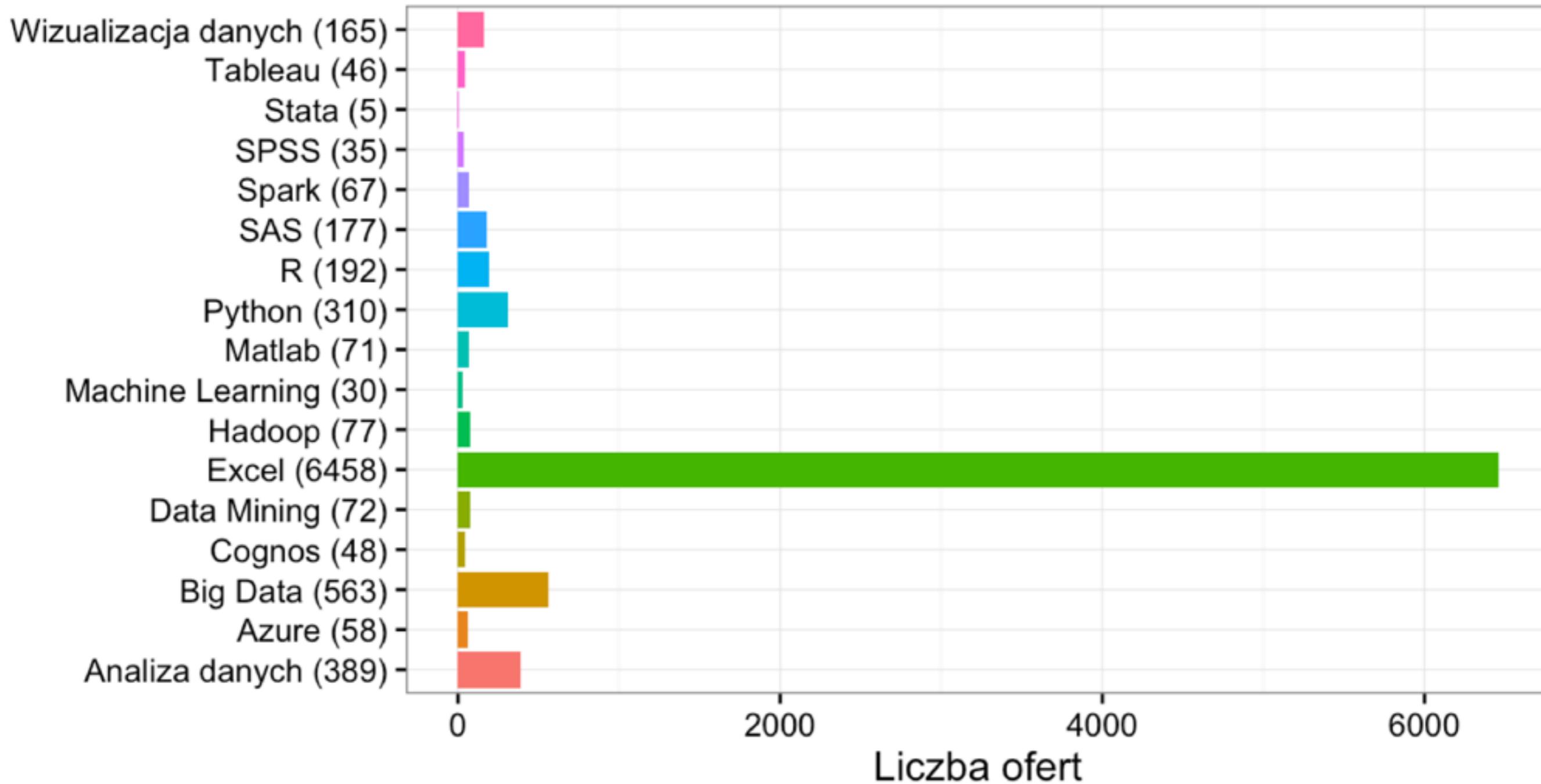
Analiza danych z pracuj.pl na 1000 ofert dla analityków:

- 326 jest z Warszawy,
- 112 z Krakowa,
- 100 z Wrocławia,
- 62 z Poznania,
- około 35 z Gdańska i Katowic.



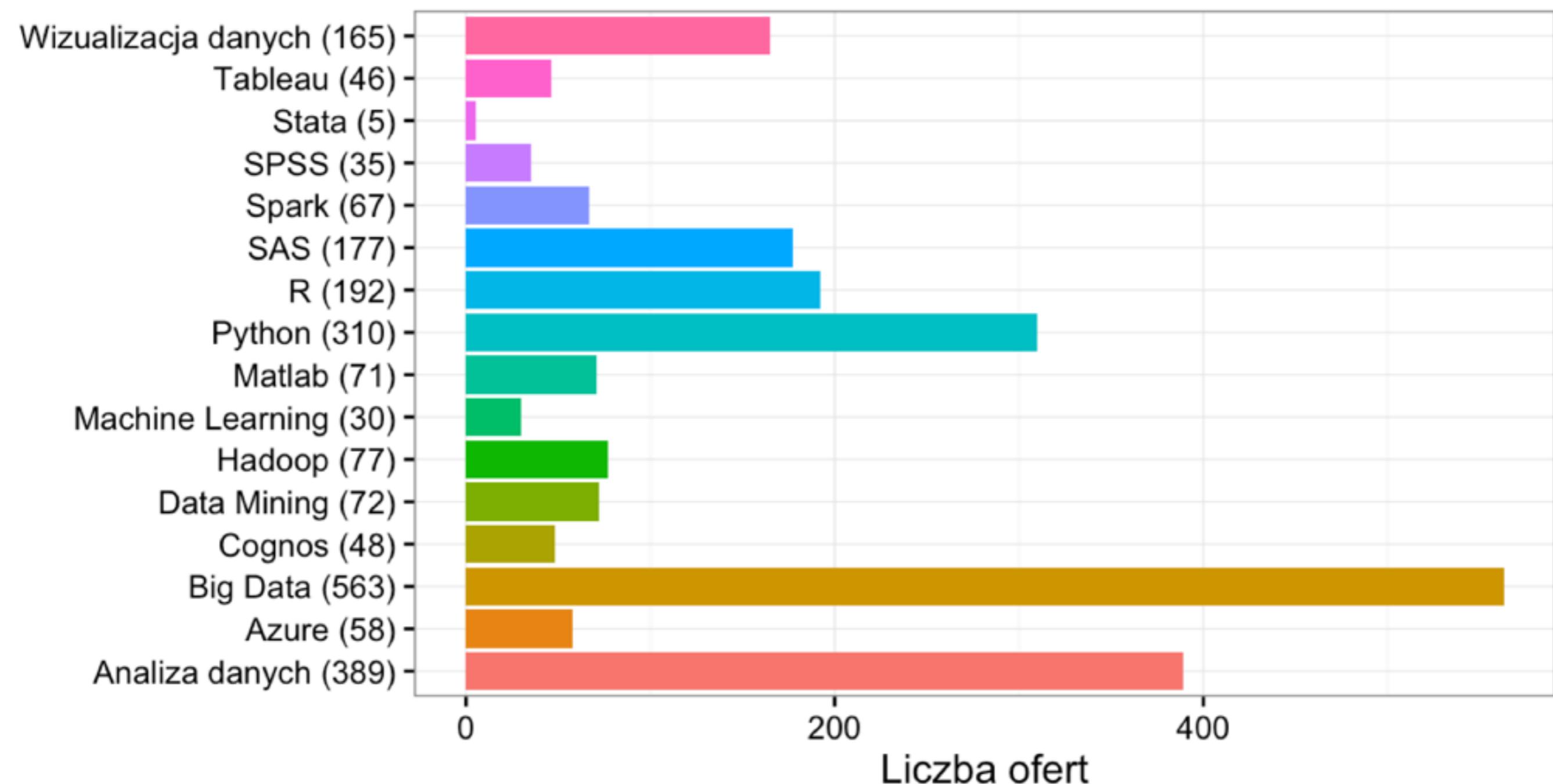
Analiza danych z pracuj.pl

Excel jest bardzo popularny w ofertach,
choć w wielu przypadkach nie dotyczy głębszych analiz



Po usunięciu Excela:

- nieznacznie więcej ofert dla R niż dla SAS,
- Matlab, SPSS, Stata mniej popularne
- dużo ofert dla Pythona, w większości dla programistów nie analityków
- Machine Learning, Big Data, takich ofert jest już całkiem sporo



BECOME DATA HERO

Looking for a job in data world? Subscribe offers newsletter.

Email Address

Subscribe

HIRE DATA HERO

Having troubles finding one?

- Access and hire within the community of data enthusiasts
- Reach out to candidates based on their skills that match your search
- Use big data analytics power and hire from datahero.tech database

Add offer

JOB'S BOARD

MACHINE LEARNING ENGINEER

THINKAPPS

San Francisco, USA

machine-learning

Added 15 days ago, 60 opens

SENIOR DATA SCIENTIST

codewise

CODEWISE

Kraków, Poland

algorithms

spark

statistics

logistic-regression

machine-learning

Added 23 days ago, 275 opens

CeNT
CENTRE
OF NEW
TECHNOLOGIES

BIOINFORMATICS, COMPUTATIONAL GENOMICS

CENTRE OF NEW TECHNOLOGIES, UNIVERSITY OF WARSAW

Warsaw, Poland

bayesian-networks

bigdata

database

matlab

Added 26 days ago, 115 opens

codilime®
CREATING VALUE

DATA SCIENTIST

CODILIME SP. Z O.O.

Warszawa

machine-learning

scikit

r

nlp

deep-learning

Added 34 days ago, 124 opens

Więcej informacji:

Odkrywać! Ujawniać! Objaśniać!

<http://biecek.pl/Eseje>

Information is Beautiful

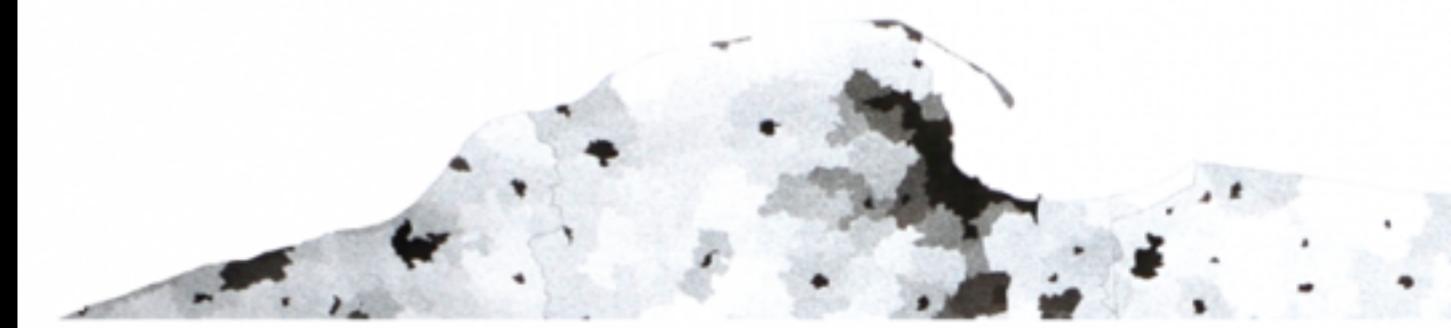
<http://www.informationisbeautiful.net/>

FlowingData

<http://flowingdata.com/>

Charts and Things

<http://chartsnthings.tumblr.com/>



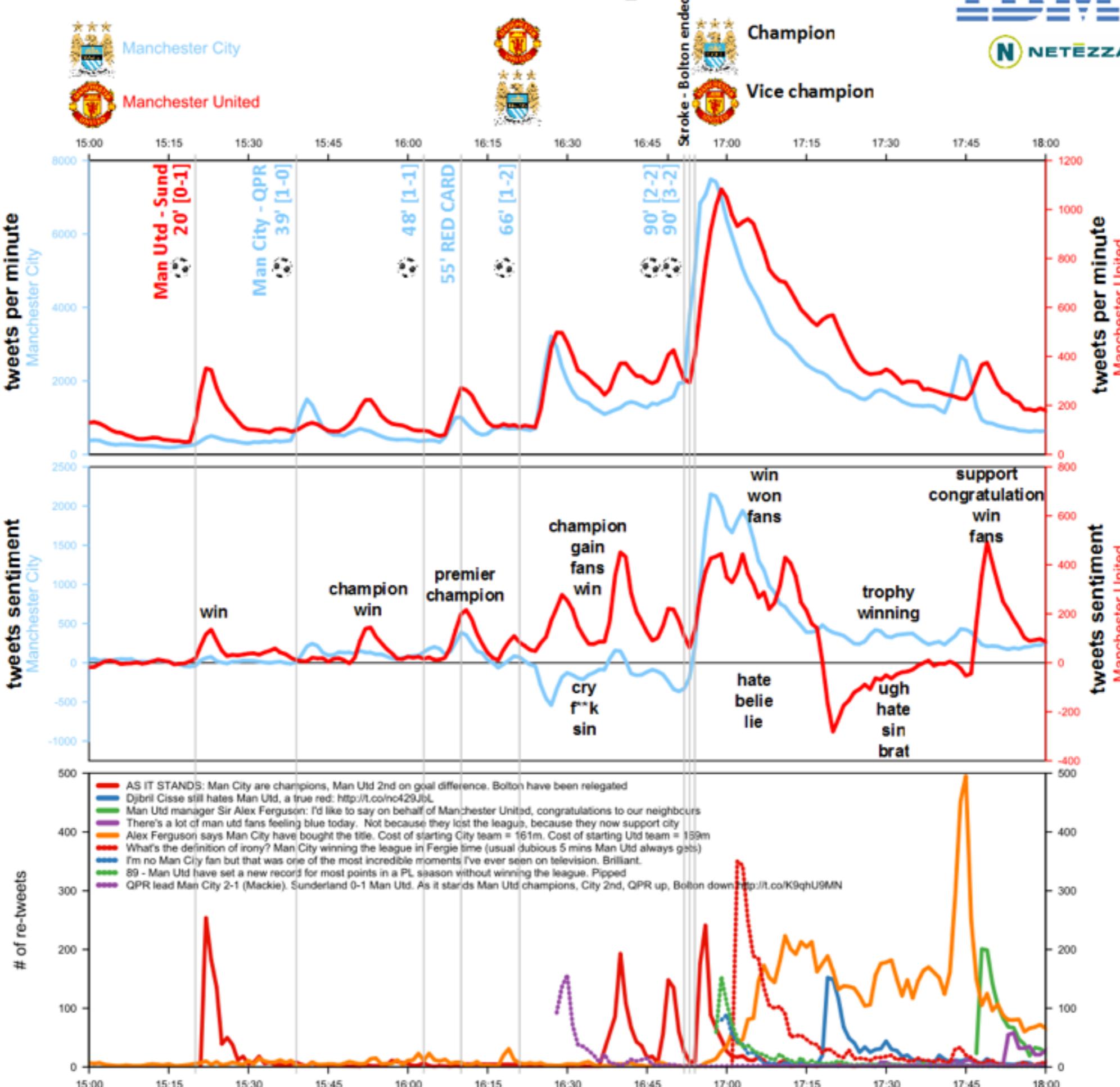
Przemysław Biecek

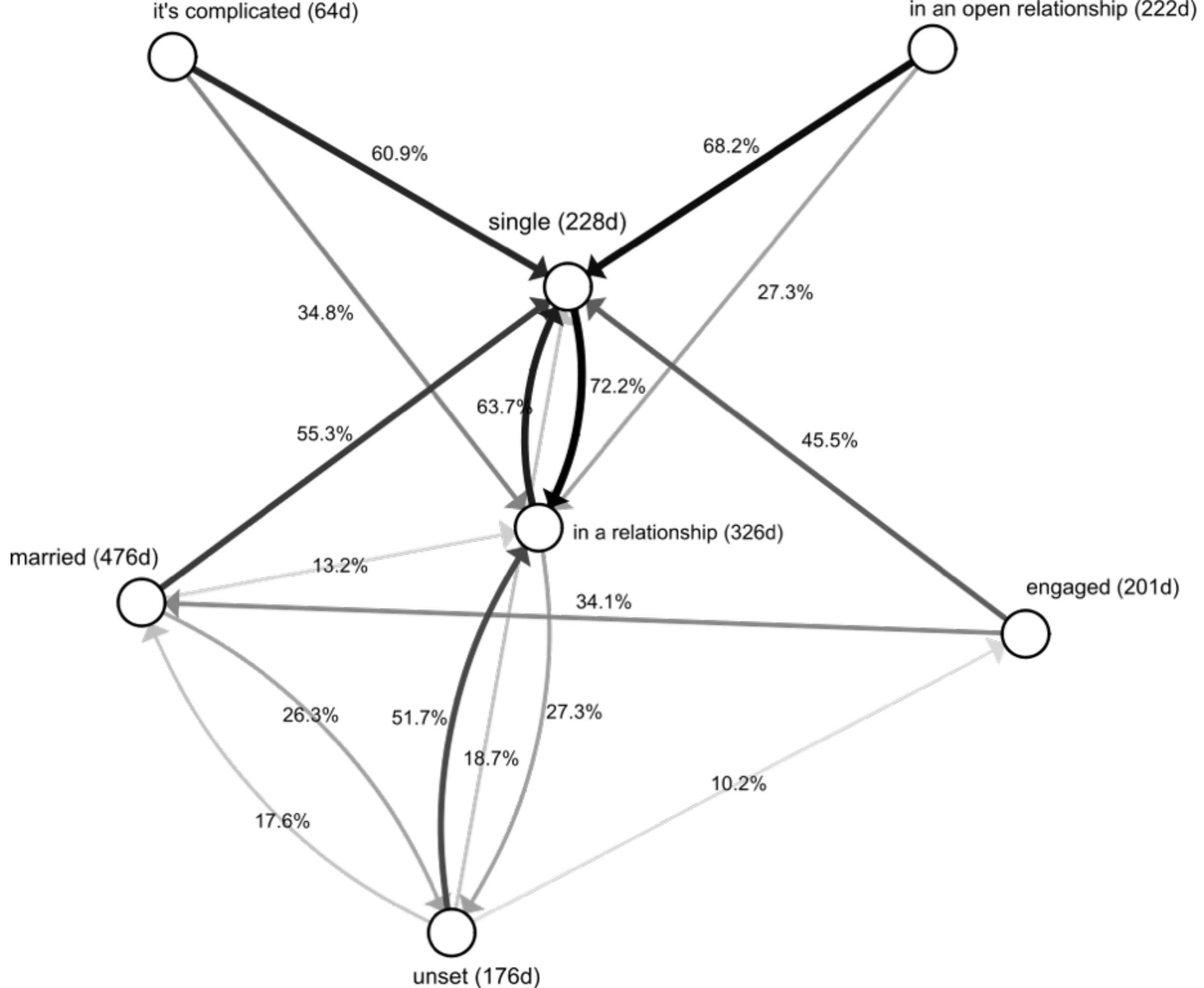
Odkrywać! Ujawniać! Objaśniać!

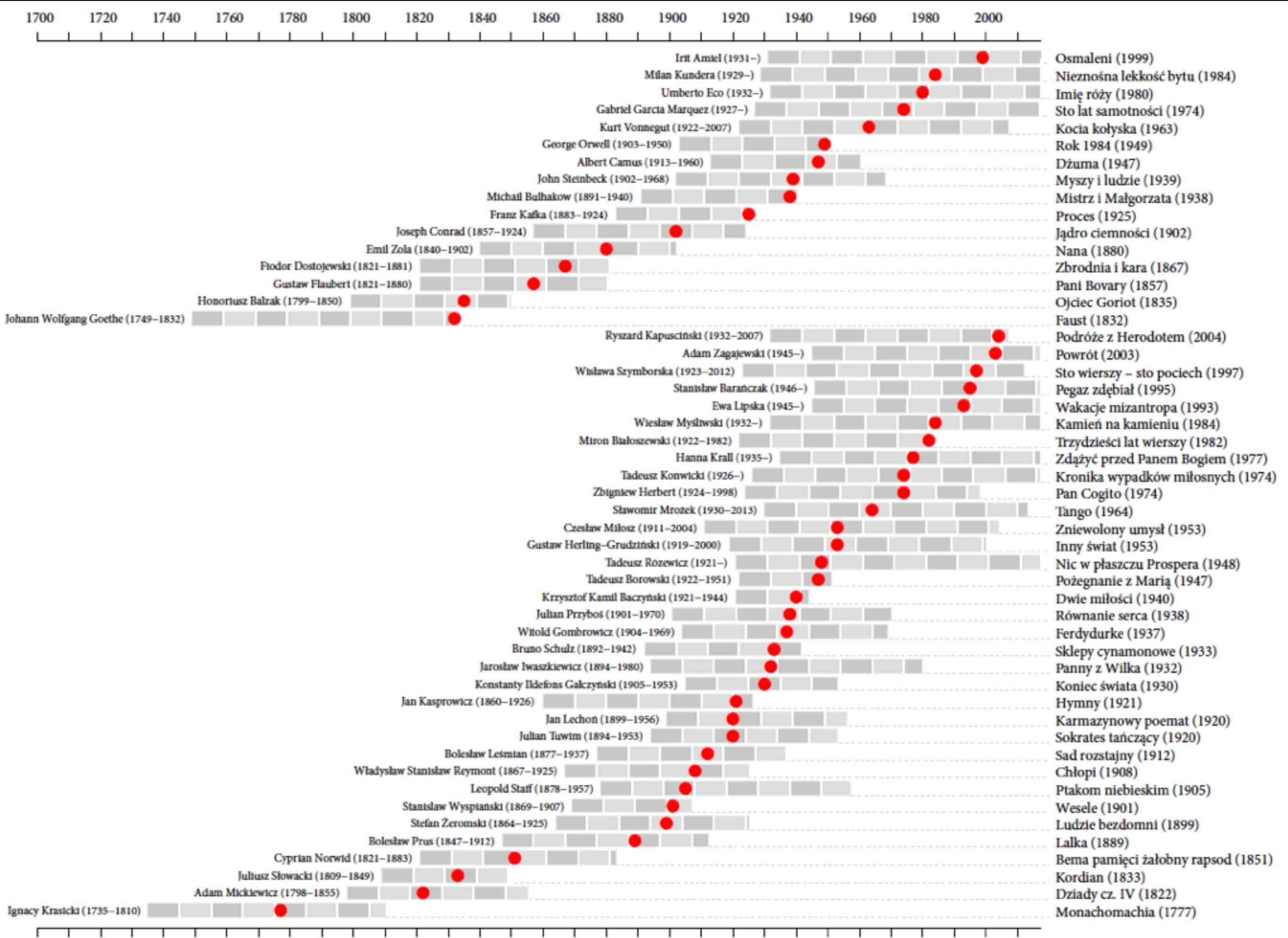
Zbiór esejów o sztuce prezentowania danych



Premier Emotions League







Pytania?

