

THE STATISTICS COMMUNICATION
SECTION PRESENTS

DATA-VIZ

A WORKSHOP SERIES ON DATA VISUALIZATION

YAN HOLTZ



CHOOSING THE RIGHT CHART FOR YOUR DATA

November 24th: 11.00 – 13.00

1

DATA VISUALIZATION IN R (BEGINNER)

December 1st: 11.00 – 13.00

2

DATA VISUALIZATION IN R (ADVANCED)

December 8th: 11.00 – 13.00

3

CEDRIC SCHERER



REGISTER AT: <https://bit.ly/3GFwgCj>

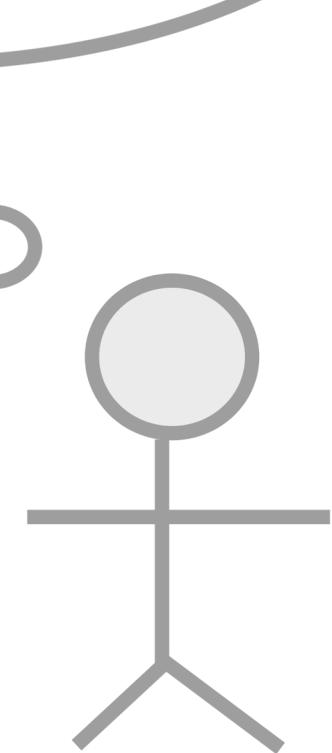
ORGANIZED BY

from Data to Viz

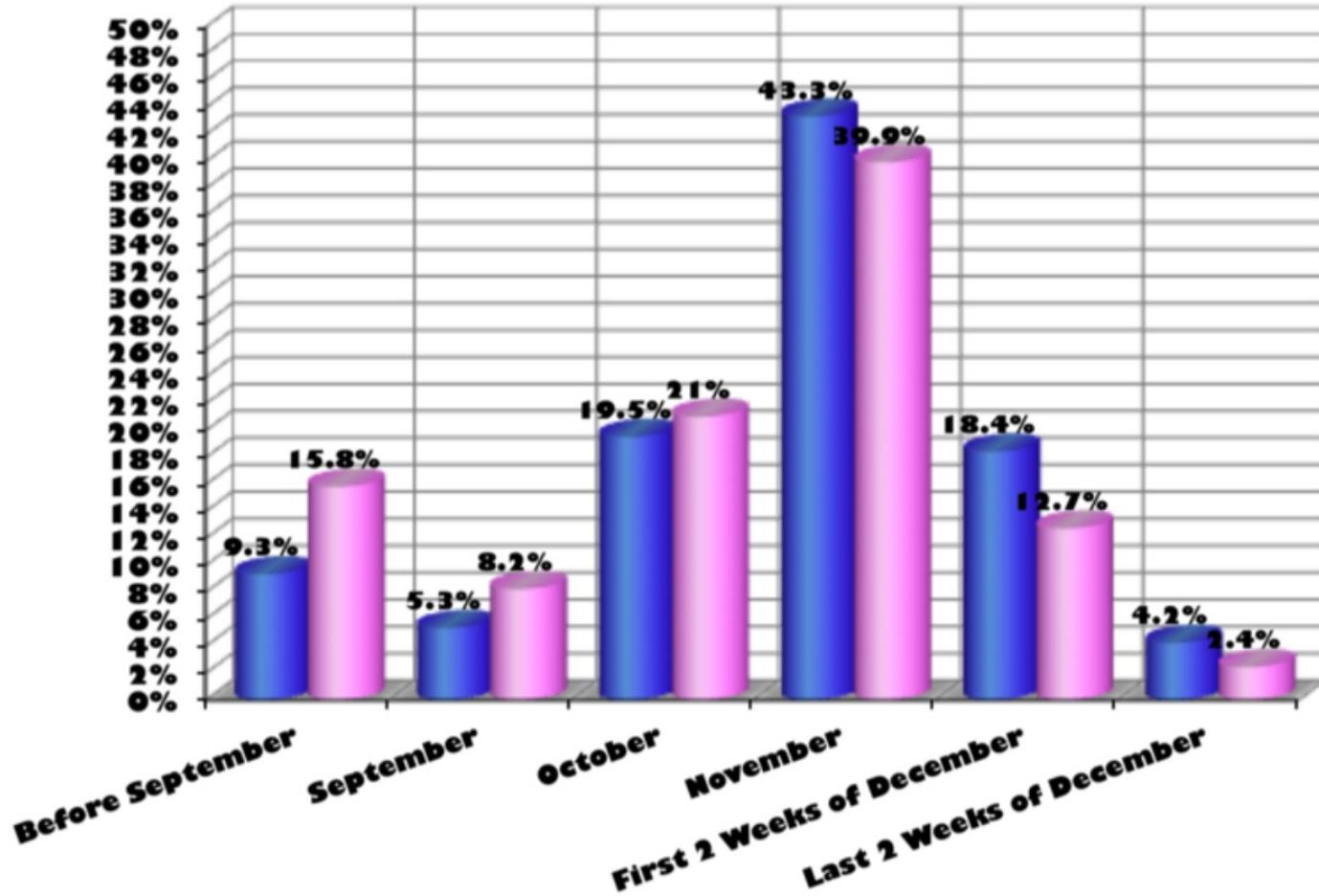
Choosing the right chart for your data

What should I do with my
data ??

Id	feature 1
A	10
B	12
C	15
...	...

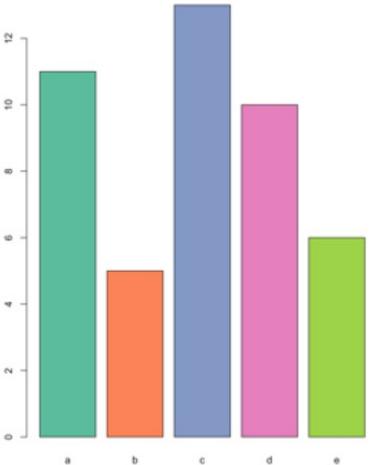


■ Men ■ Women



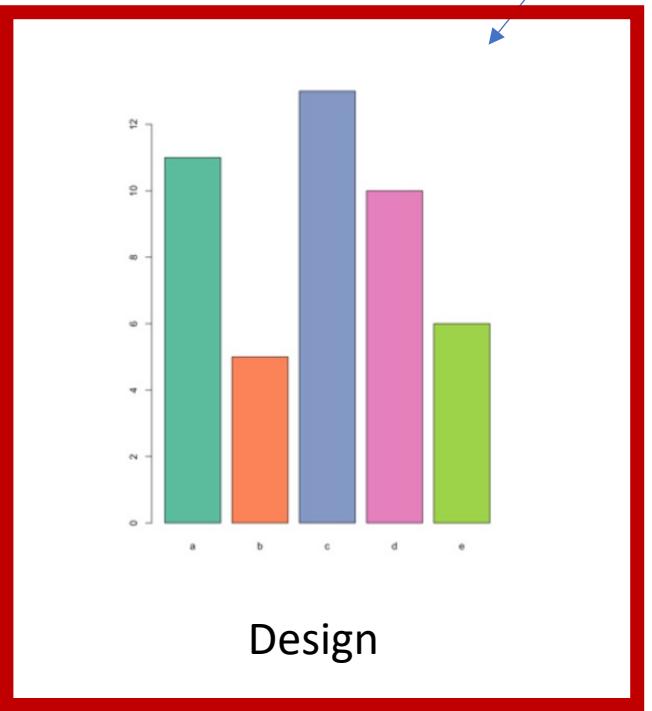
...

What about
that?



Design

Code



A chart, why?

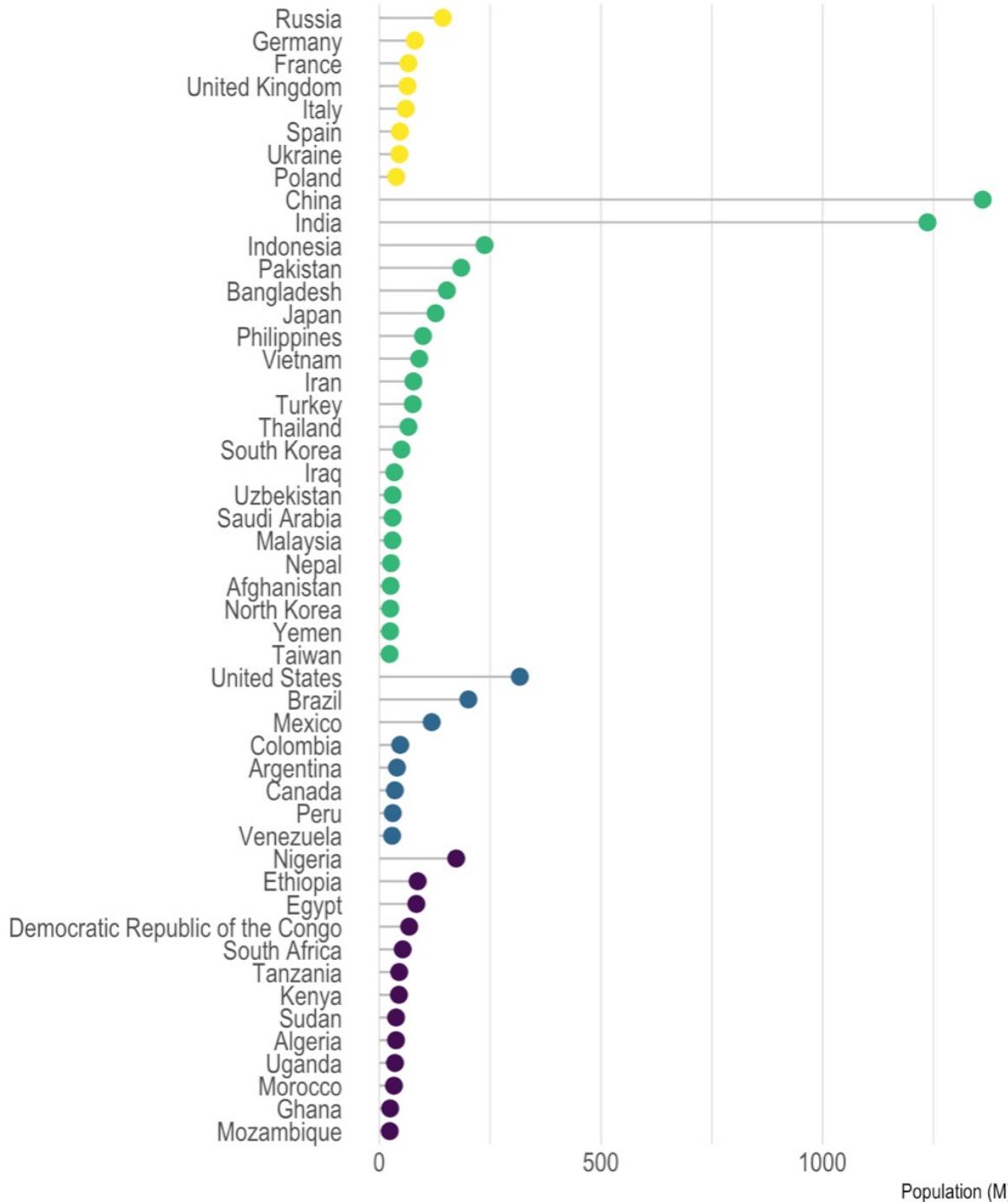
World population
distribution

	A	B	C	D
1	region	subregion	key	value
2	Asia	Southern Asia	Afghanistan	25500100
3	Europe	Northern Europe	Åland Islands	28502
4	Europe	Southern Europe	Albania	2821977
5	Africa	Northern Africa	Algeria	37900000
6	Oceania	Polynesia	American Samoa	55519
7	Europe	Southern Europe	Andorra	76246
8	Africa	Middle Africa	Angola	20609294
9	Americas	Caribbean	Anguilla	13452
10			Antarctica	-1
11	Americas	Caribbean	Antigua and Barbuda	86295
12	Americas	South America	Argentina	40117096
13	Asia	Western Asia	Armenia	3024100
14	Americas	Caribbean	Aruba	101484
15	Oceania	Australia and New Zealand	Australia	23254142
16	Europe	Western Europe	Austria	8501502
17	Asia	Western Asia	Azerbaijan	9235100
18	Americas	Caribbean	Bahamas	-1
19	Asia	Western Asia	Bahrain	1234571
20	Asia	Southern Asia	Bangladesh	152518015
21	Americas	Caribbean	Barbados	274200
22	Europe	Eastern Europe	Belarus	9465500
23	Europe	Western Europe	Belgium	11175653
24	Americas	Central America	Belize	312971
25	Africa	Western Africa	Benin	10323000
26	Americas	Northern America	Bermuda	64237
27	Asia	Southern Asia	Bhutan	740990
28	Americas	South America	Bolivia	10027254
29	Americas	Caribbean	Bonaire	-1
30	Europe	Southern Europe	Bosnia and Herzegovina	3791622
31	Africa	Southern Africa	Botswana	2024904
32			Bouvet Island	-1
33	Americas	South America	Brazil	201032714
34	Africa	Eastern Africa	British Indian Ocean Territory	-1
35	Americas	Caribbean	British Virgin Islands	29537

- Biggest country?
- Rank?
- Distribution by continent?



“Picture superiority effect”

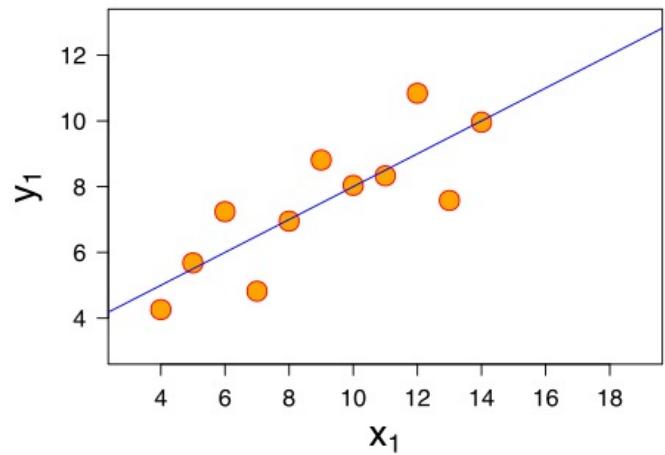


$$Y = 3 + 0.5x$$
$$\text{Cor} = 0.8$$

OK, but we have
summary statistics!

$$\text{Mean}(x) = 9$$
$$\text{Var}(x) = 11$$

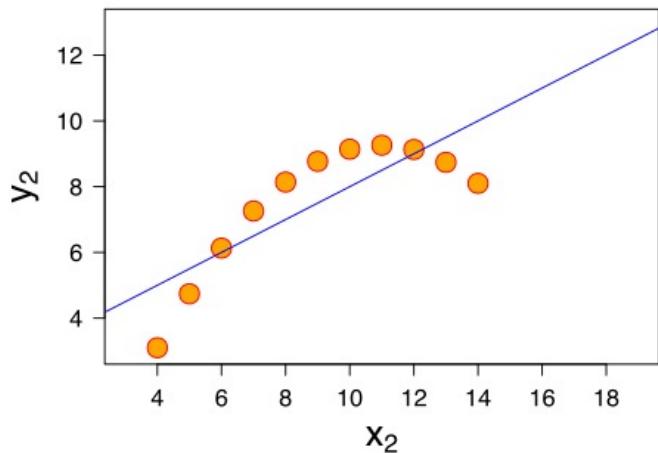
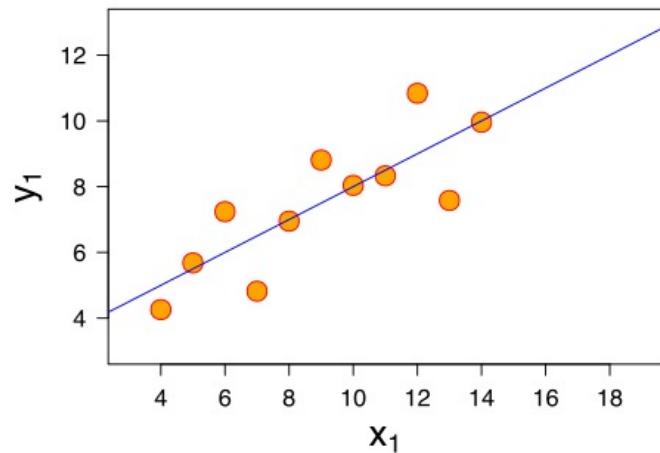
$$\text{Mean}(Y) = 7.5$$
$$\text{Var}(Y) = 4.1$$



$$Y = 3 + 0.5x$$
$$\text{Cor} = 0.8$$

$$\text{Mean}(x) = 9$$
$$\text{Var}(x) = 11$$

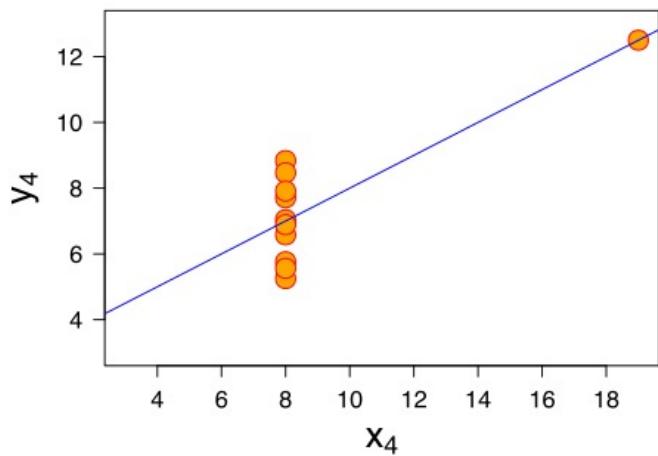
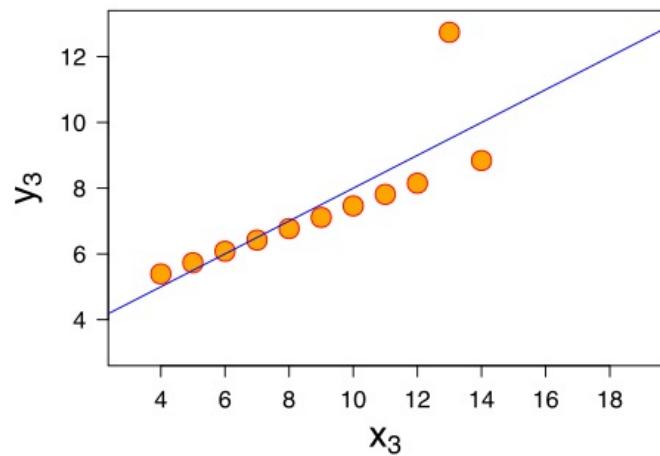
$$\text{Mean}(Y) = 7.5$$
$$\text{Var}(Y) = 4.1$$



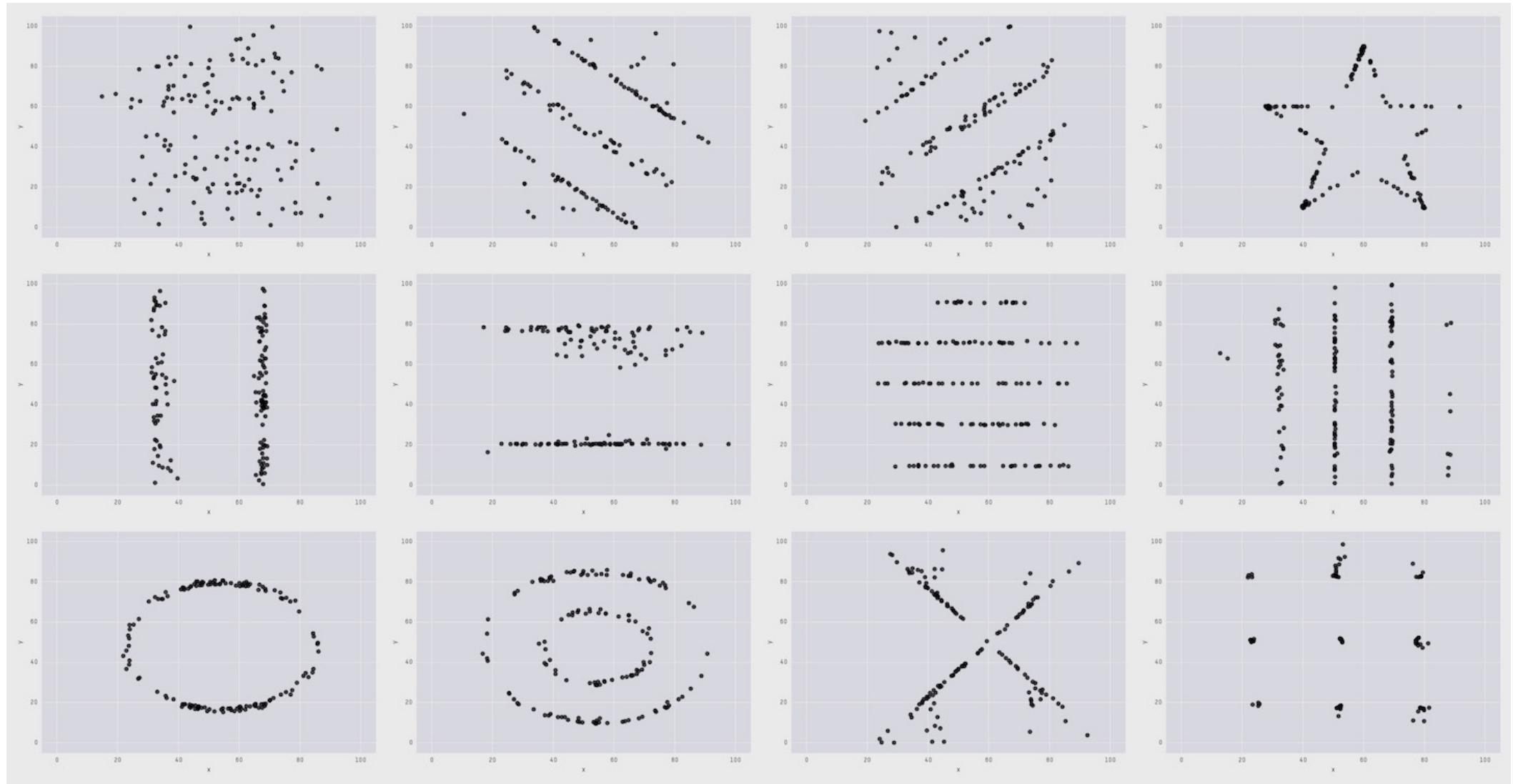
$Y = 3 + 0.5x$
Cor = 0.8

Mean(x) = 9
Var(x) = 11

Mean(Y) = 7.5
Var(Y) = 4.1

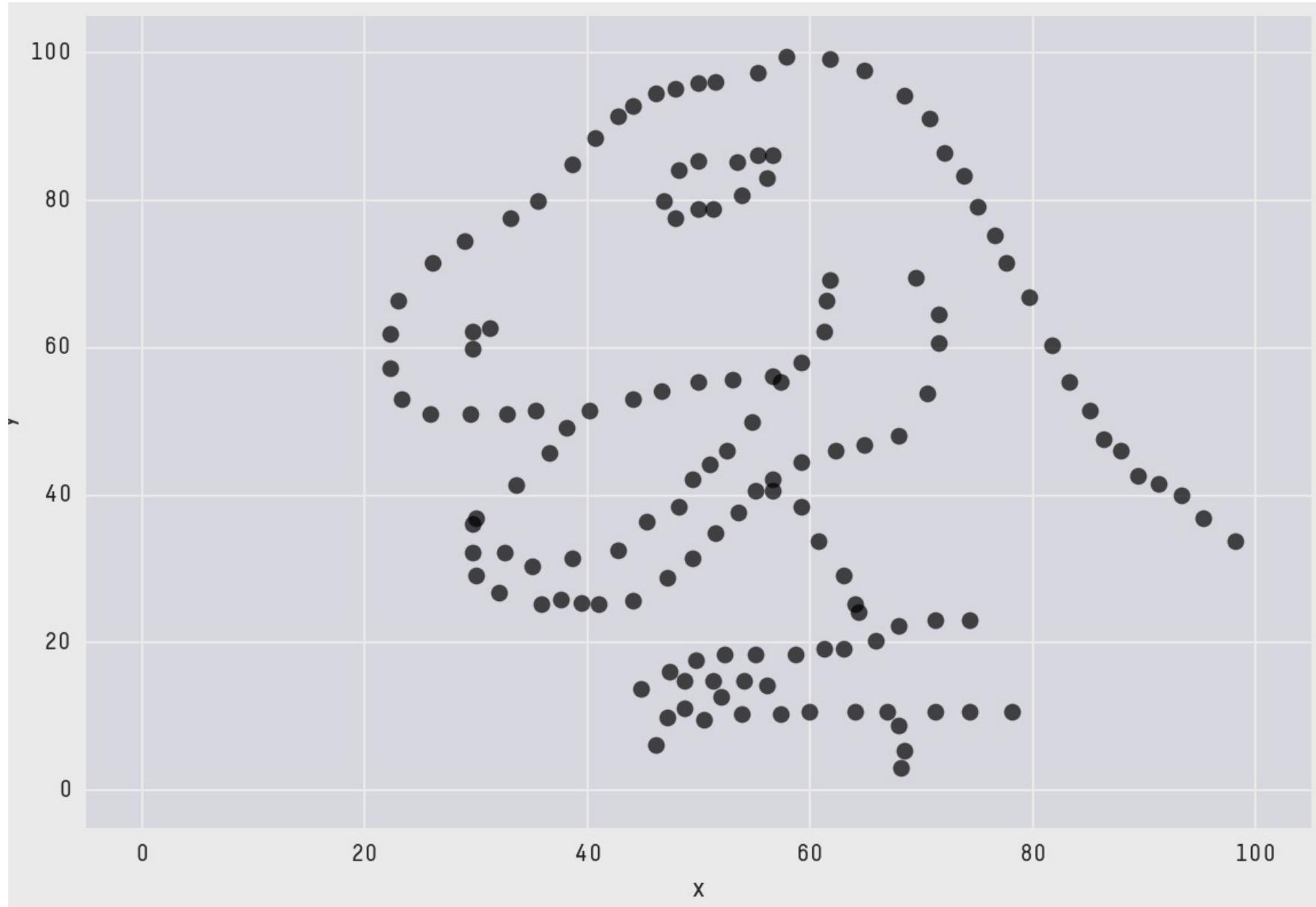


Anscombe's quartet



The Datasaurus Dozen

Cairo, 2017
Matejka & Fitzmaurice, 2017



The Datasaurus Dozen

Cairo, 2017

Matejka & Fitzmaurice, 2017

Looking for a chart ?

What you
can do

What you
should do

Caveats to
avoid

How to
build it

WHAT YOU CAN DO

A classification of chart types based on data input format



Scatterplot



Scatterplot



2d density chart

Who sells more weapons ?

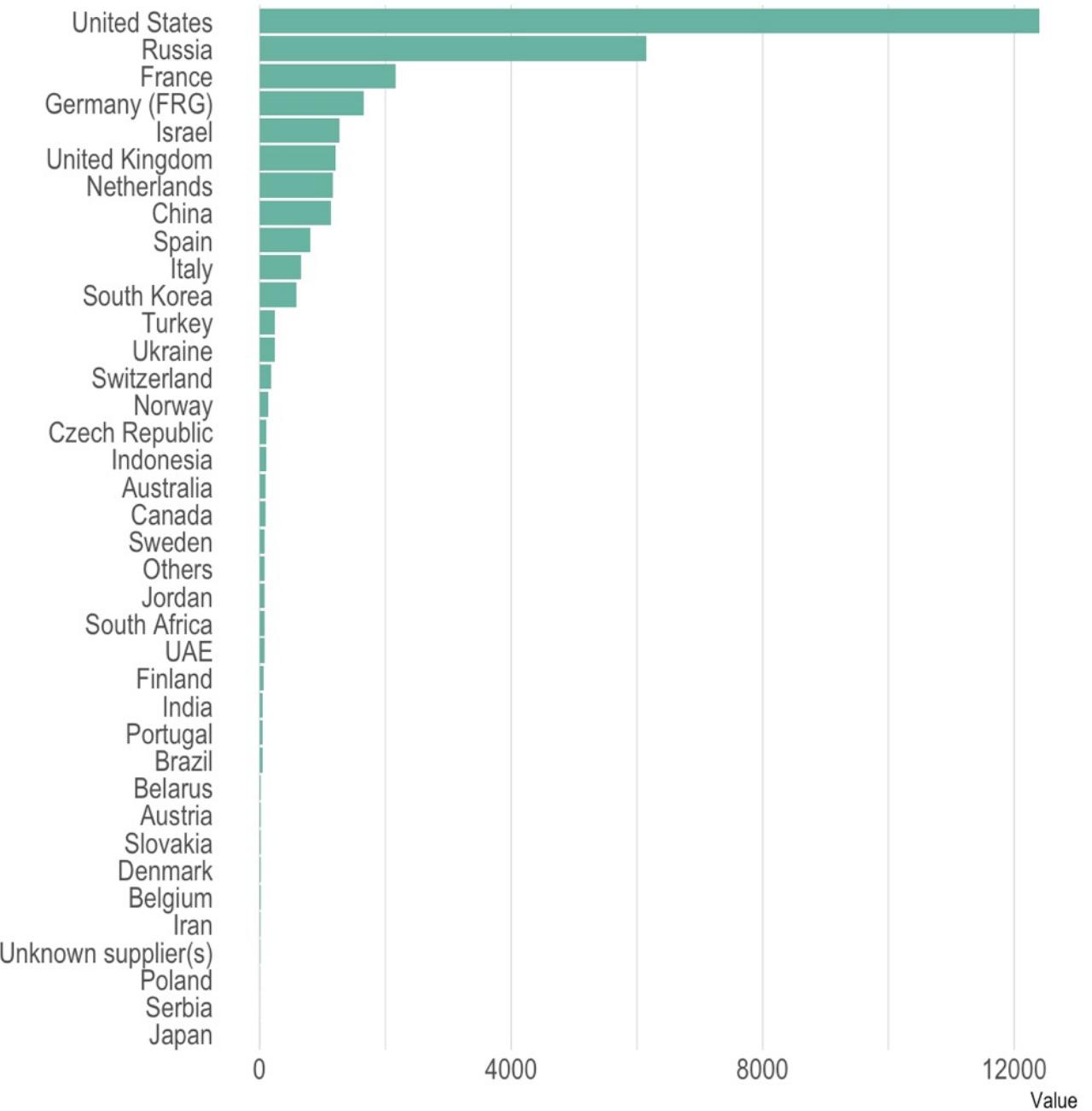
Country	Value
United States	12394
Russia	6148
Germany (FRG)	1653
France	2162
United Kingdom	1214
China	1131



?

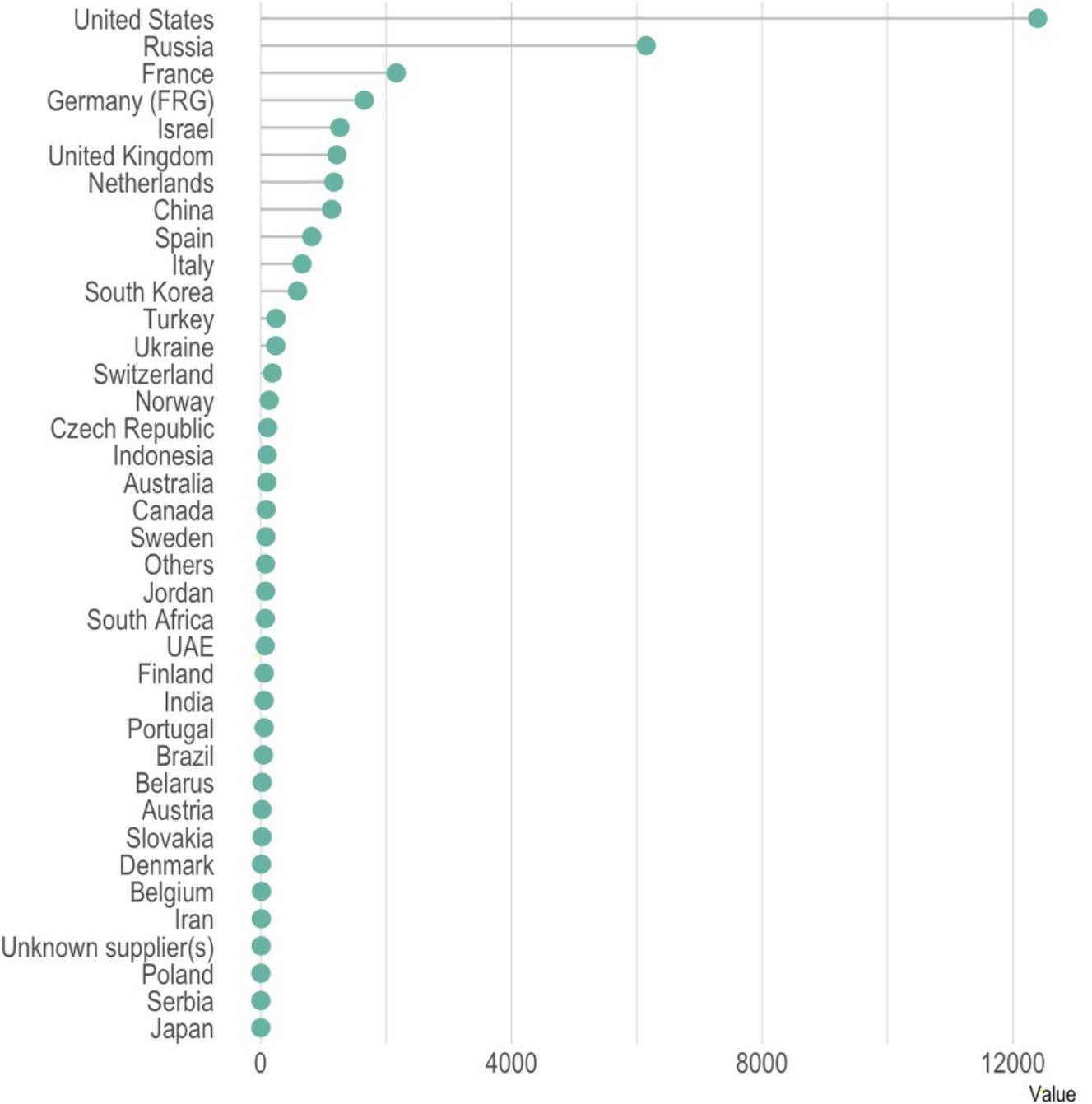
Who sells more weapons ?

Country	Value
United States	12394
Russia	6148
Germany (FRG)	1653
France	2162
United Kingdom	1214
China	1131



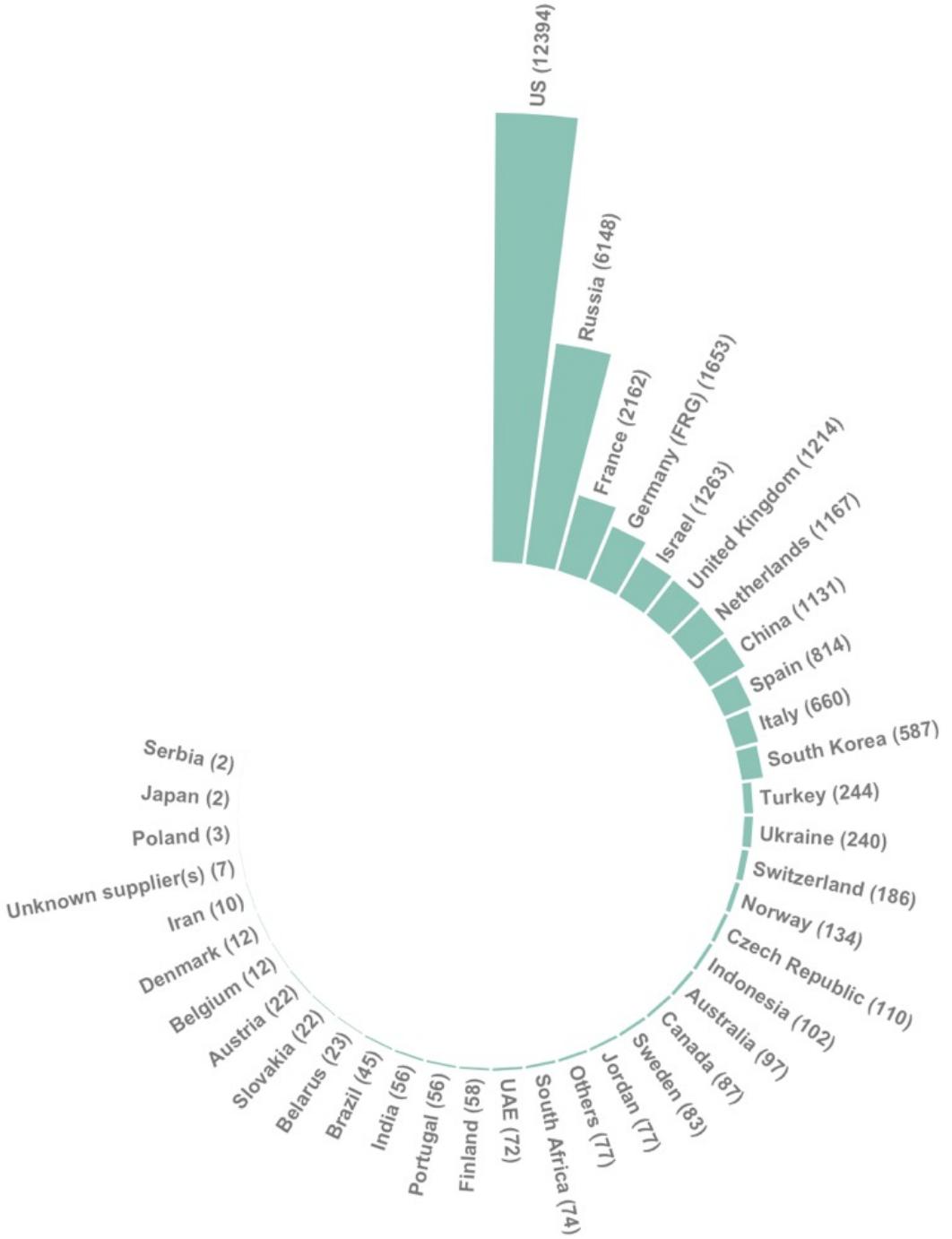
Who sells more weapons ?

Country	Value
United States	12394
Russia	6148
Germany (FRG)	1653
France	2162
United Kingdom	1214
China	1131



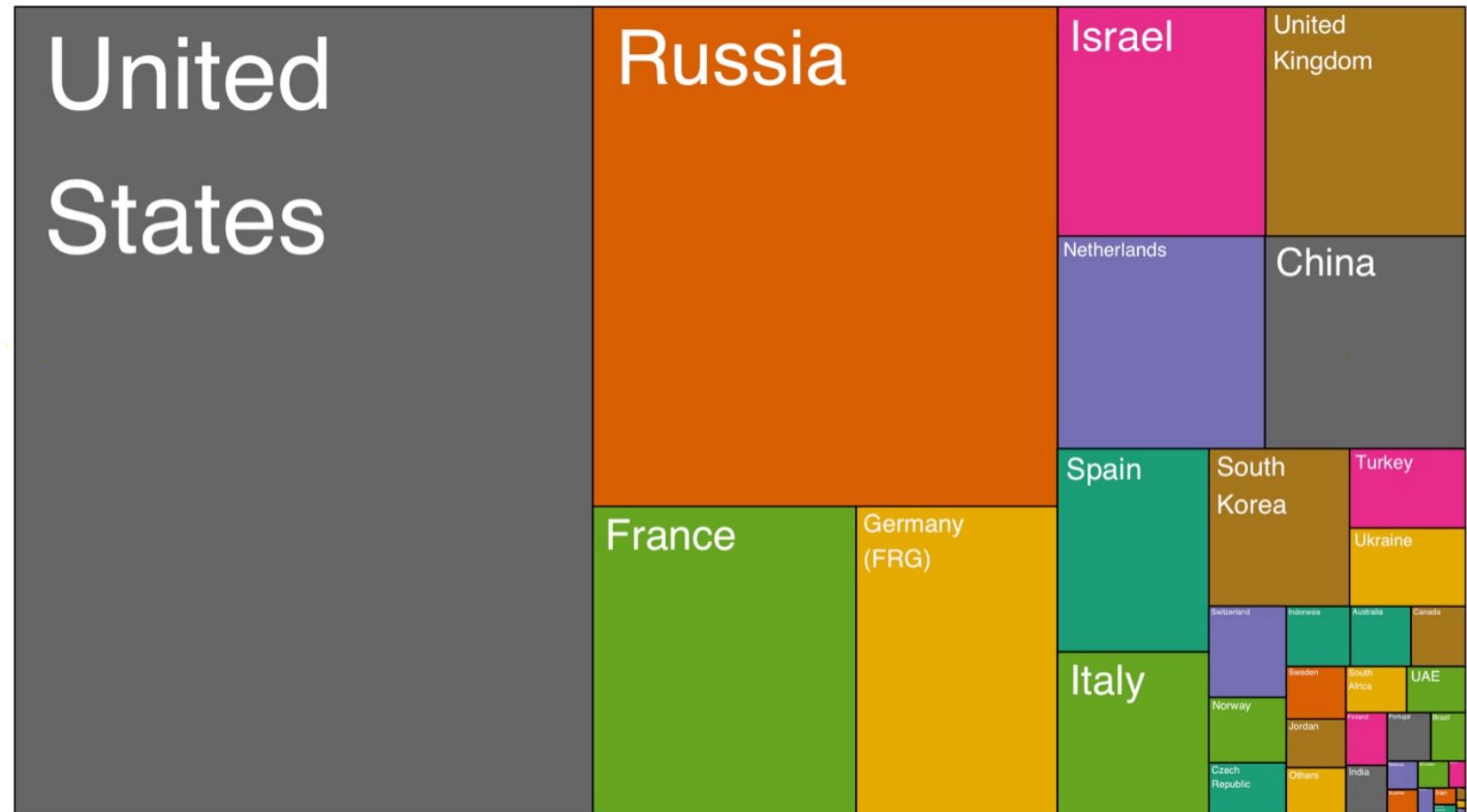
Who sells more weapons ?

Country	Value
United States	12394
Russia	6148
Germany (FRG)	1653
France	2162
United Kingdom	1214
China	1131



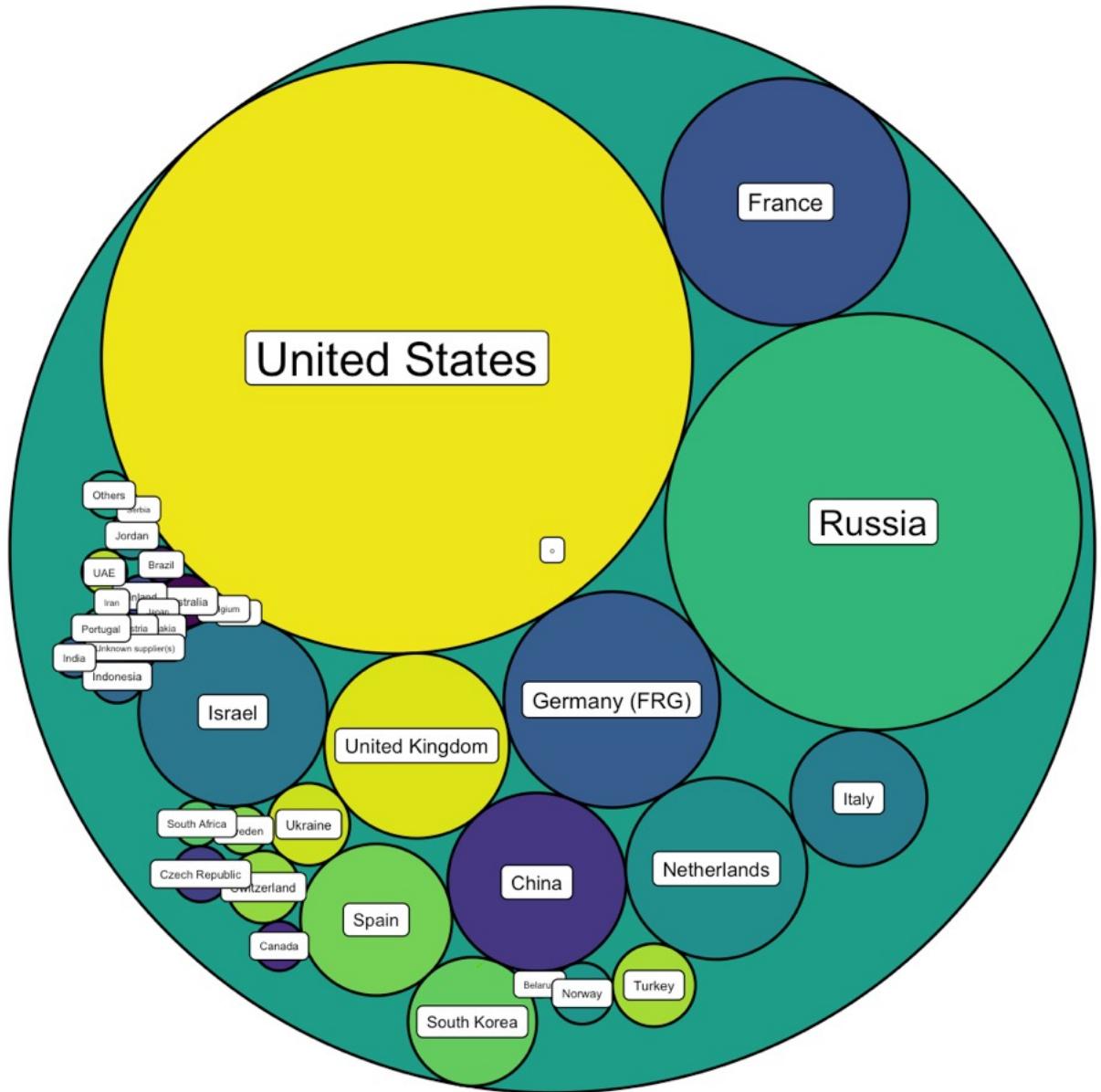
Who sells more weapons ?

Country	Value
United States	12394
Russia	6148
Germany (FRG)	1653
France	2162
United Kingdom	1214
China	1131



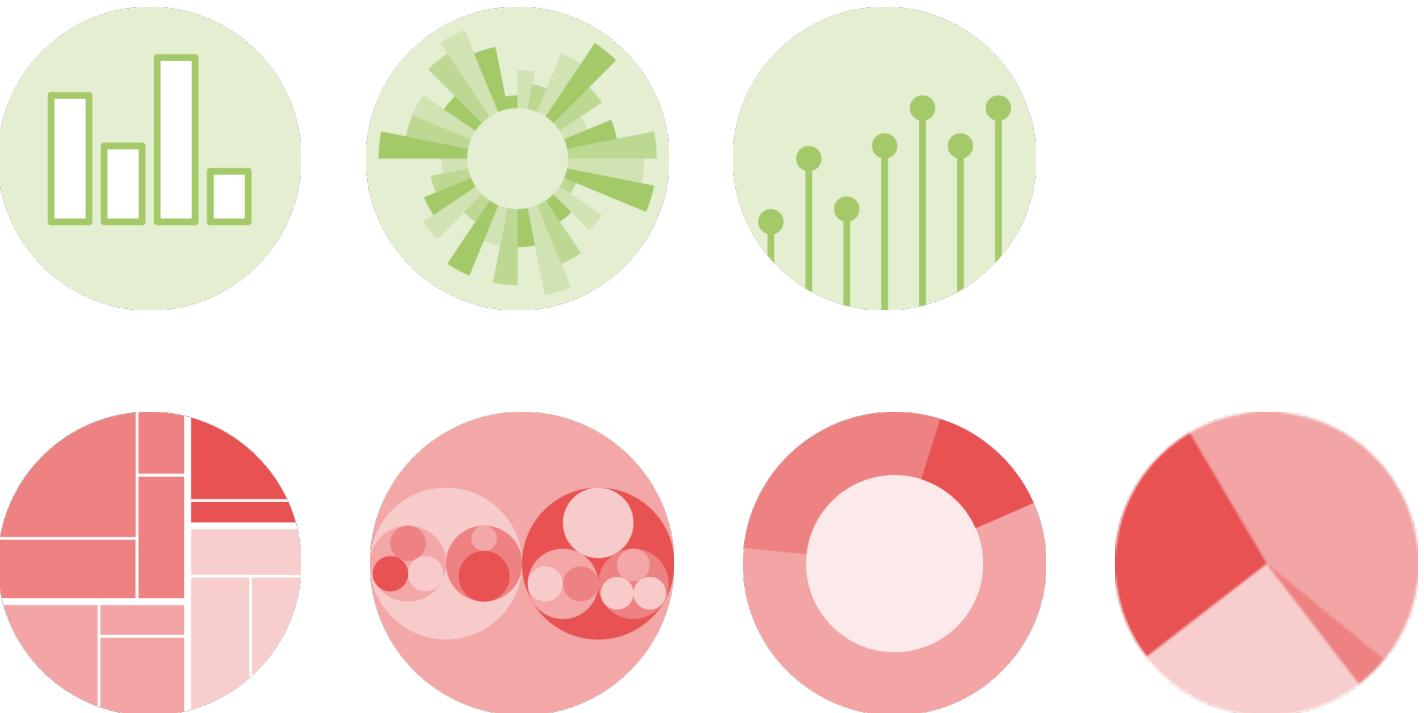
Who sells more weapons ?

Country	Value
United States	12394
Russia	6148
Germany (FRG)	1653
France	2162
United Kingdom	1214
China	1131



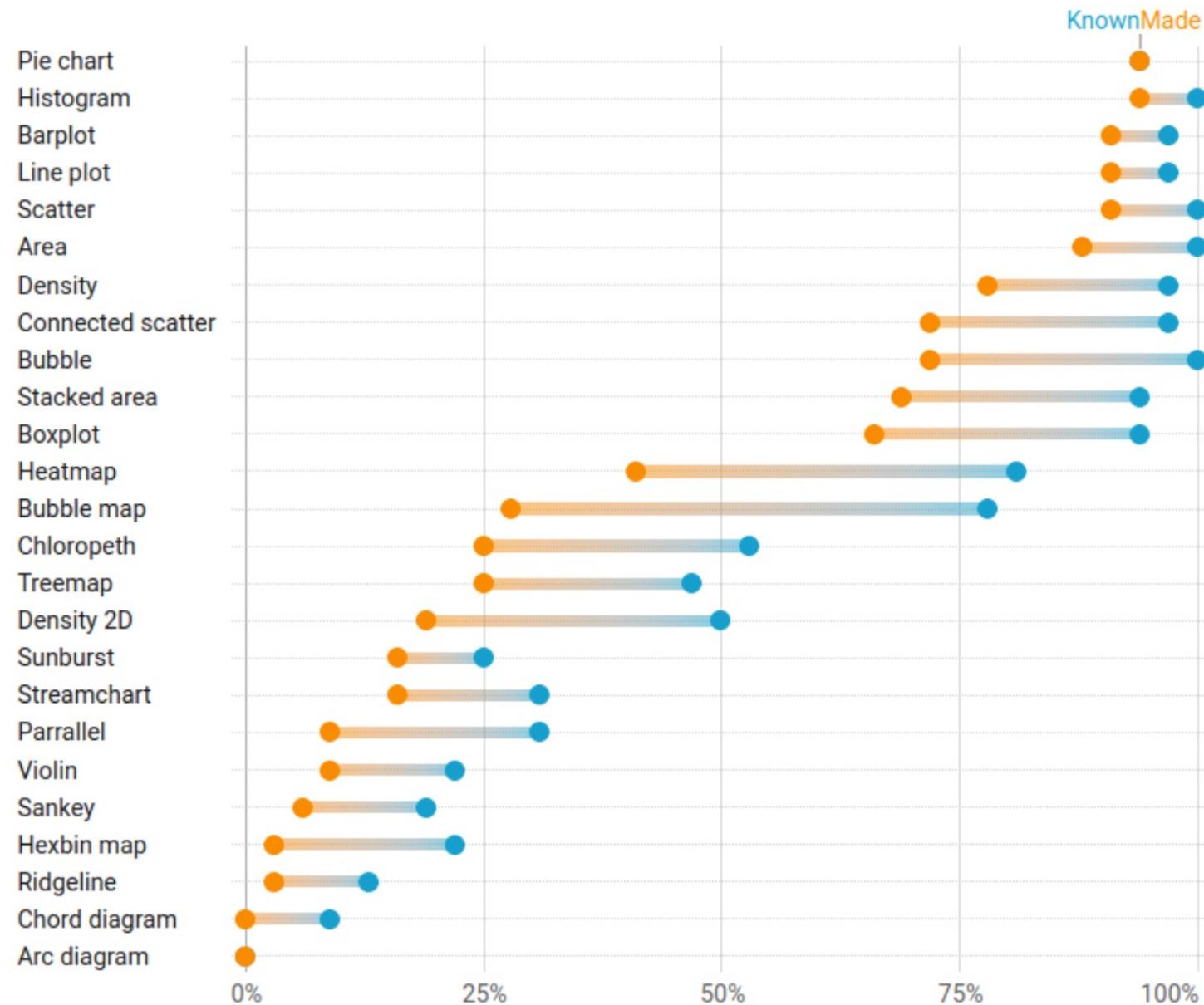
Who sells more weapons ?

Country	Value
United States	12394
Russia	6148
Germany (FRG)	1653
France	2162
United Kingdom	1214
China	1131



Charts you made and charts you know 18/12

Dataviz training session 18/12

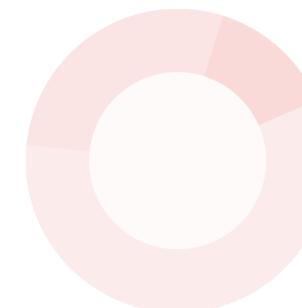
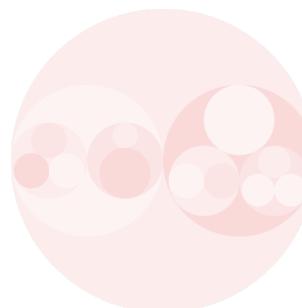
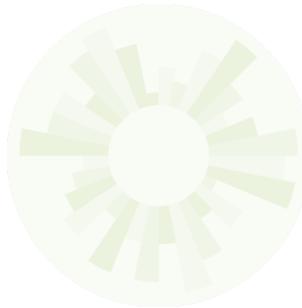


Who sells more weapons ?

Country	Value
United States	12394
Russia	6148
Germany (FRG)	1653
France	2162
United Kingdom	1214
China	1131

1 categorical variable

1 numerical variable



1 observation per group

Perception of probability

text	value
Improbable	33
Almost Certainly	98
Likely	60
Almost Certainly	98
Unlikely	10
Probably Not	25
About Even	50
Probably	75

Perception of probability

text	value
Improbable	33
Almost Certainly	98
Likely	60
Almost Certainly	98
Unlikely	10
Probably Not	25
About Even	50
Probably	75

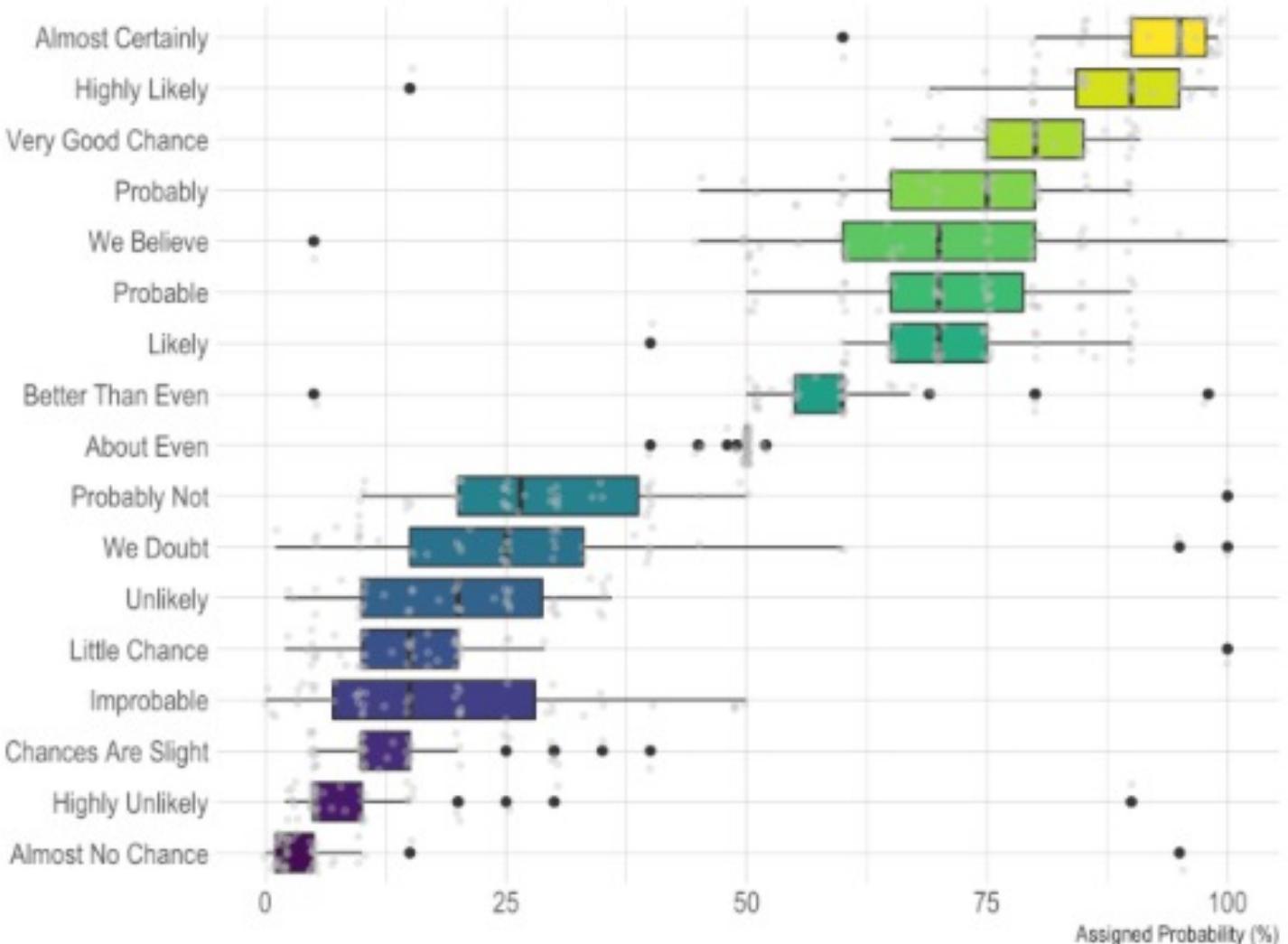
Several
observations
per group

1 categorical
variable

1 numeric
variable

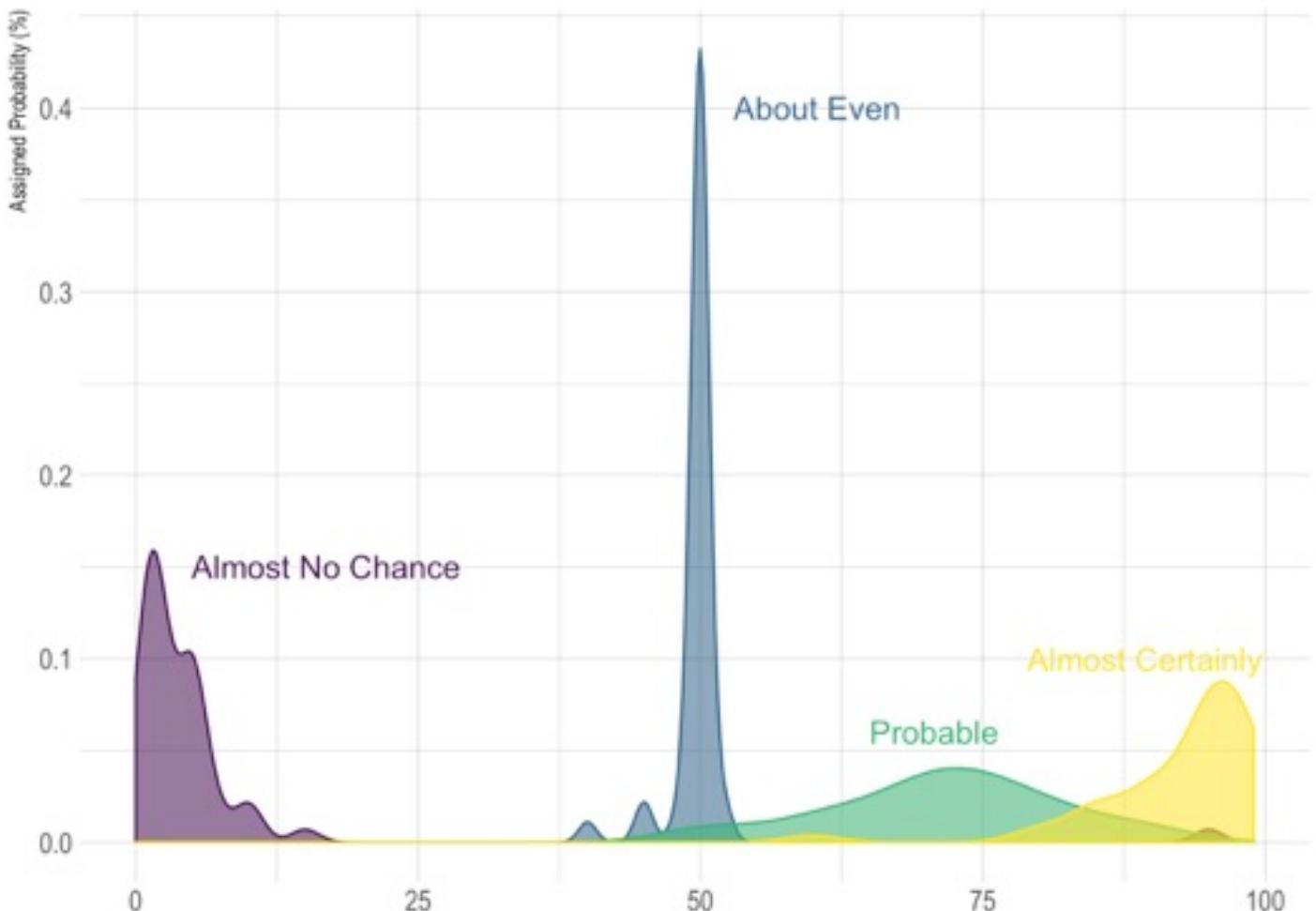
Perception of probability

text	value
Improbable	33
Almost Certainly	98
Likely	60
Almost Certainly	98
Unlikely	10
Probably Not	25
About Even	50
Probably	75



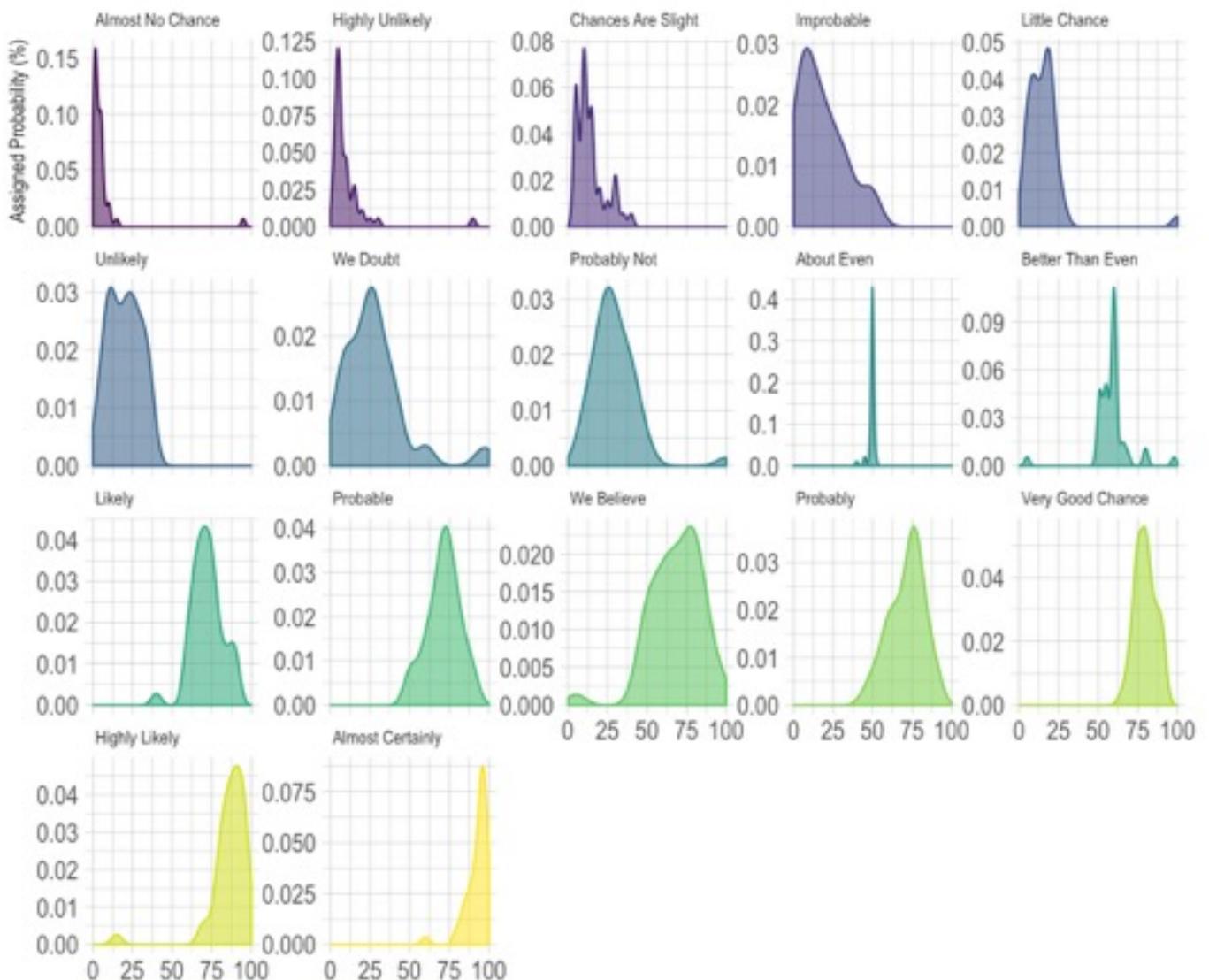
Perception of probability

text	value
Improbable	33
Almost Certainly	98
Likely	60
Almost Certainly	98
Unlikely	10
Probably Not	25
About Even	50
Probably	75



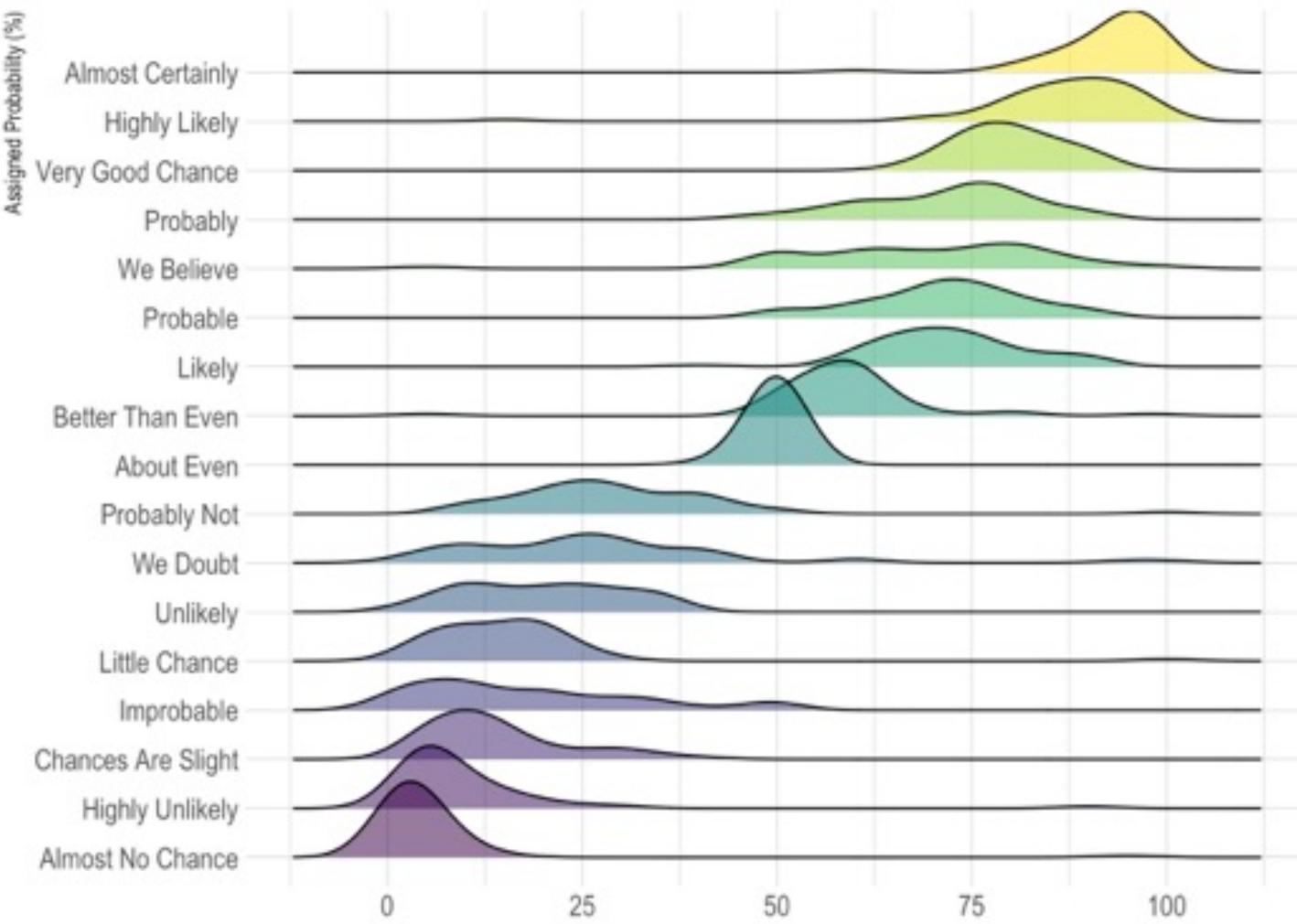
Perception of probability

text	value
Improbable	33
Almost Certainly	98
Likely	60
Almost Certainly	98
Unlikely	10
Probably Not	25
About Even	50
Probably	75



Perception of probability

text	value
Improbable	33
Almost Certainly	98
Likely	60
Almost Certainly	98
Unlikely	10
Probably Not	25
About Even	50
Probably	75



So...

- Knowing the **possibilities** is the **first step** in chart choice
- Not easy to know **all** the **chart types**
- Hard to figure out options **from a dataset**



Let's build a decision tree



from Data to Viz

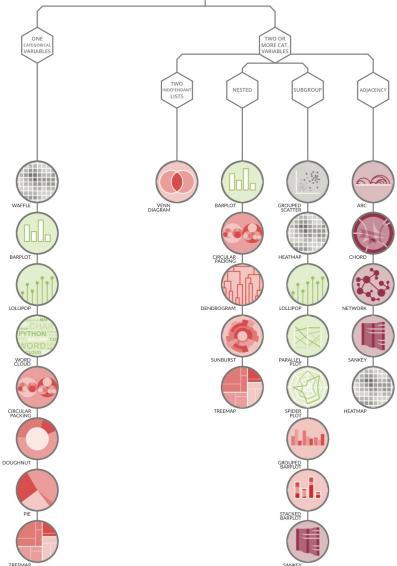
'From Data to Viz' is a classification of chart types based on input data format. It will help you find the perfect chart in three simple steps :

- 1 Identify what type of data you have.
 - 2 Go to the corresponding decision tree and follow it down to a set of possible charts.
 - 3 Choose the chart from the set that will suit your data and your needs best.

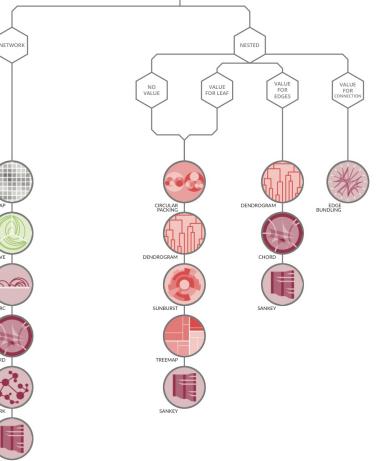
Dataviz is a world with endless possibilities and this project does not claim to be exhaustive. However it should provide you with a good starting point. For an interactive version and much more, visit:

data-to-viz.com

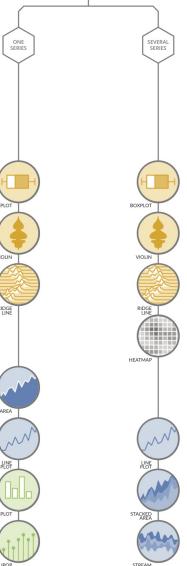
CATEGORIC



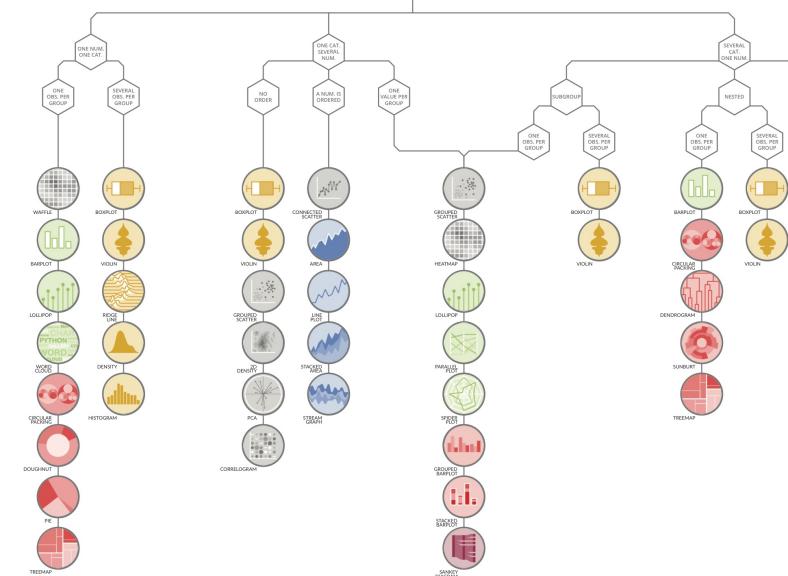
RELATIONAL



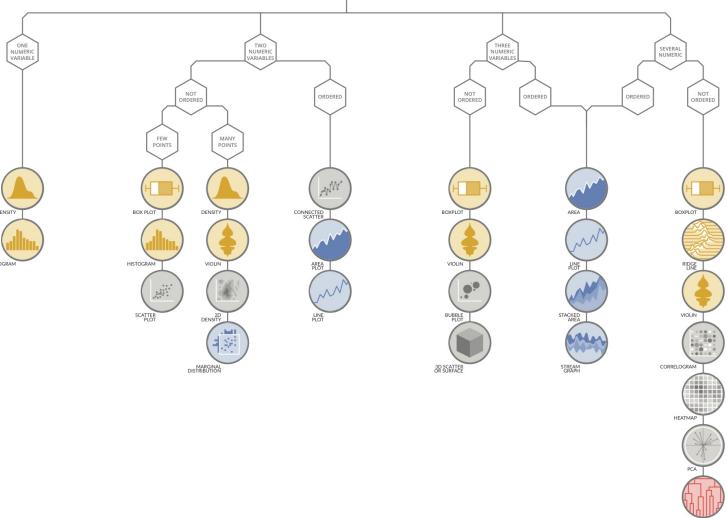
TIME SERIES



CATEGORIC AND NUMERIC



NUMERIC



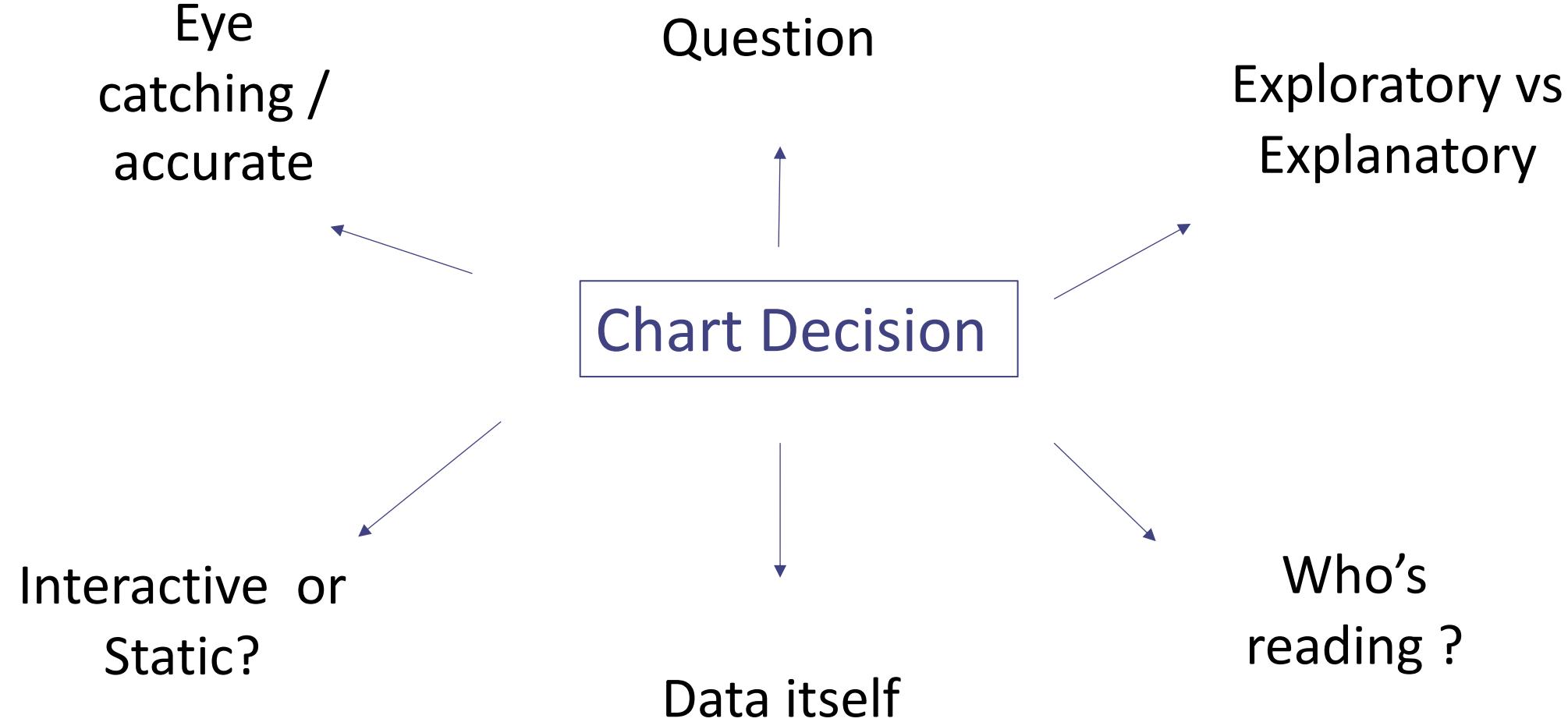
Data-to-viz.com

There are no limits in dataviz!

<https://xeno.graphics/>

WHAT YOU SHOULD DO

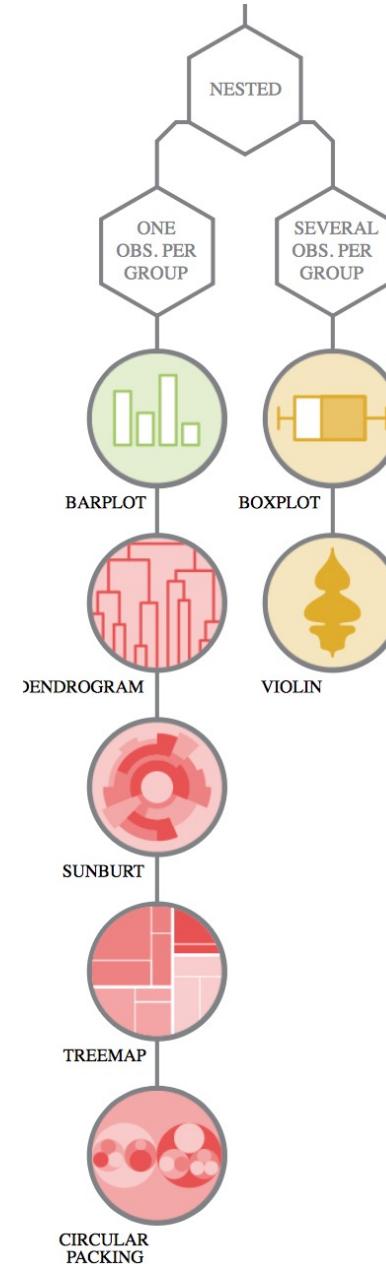
About 20 examples of storytelling with data



Question – Explo/Expla – Reader – Data – Interactivity – Eye catching

WHAT DO YOU WANT TO SHOW ?

- Distribution
- Evolution
- Correlation
- Maps
- Ranking
- Flow
- Part of a whole

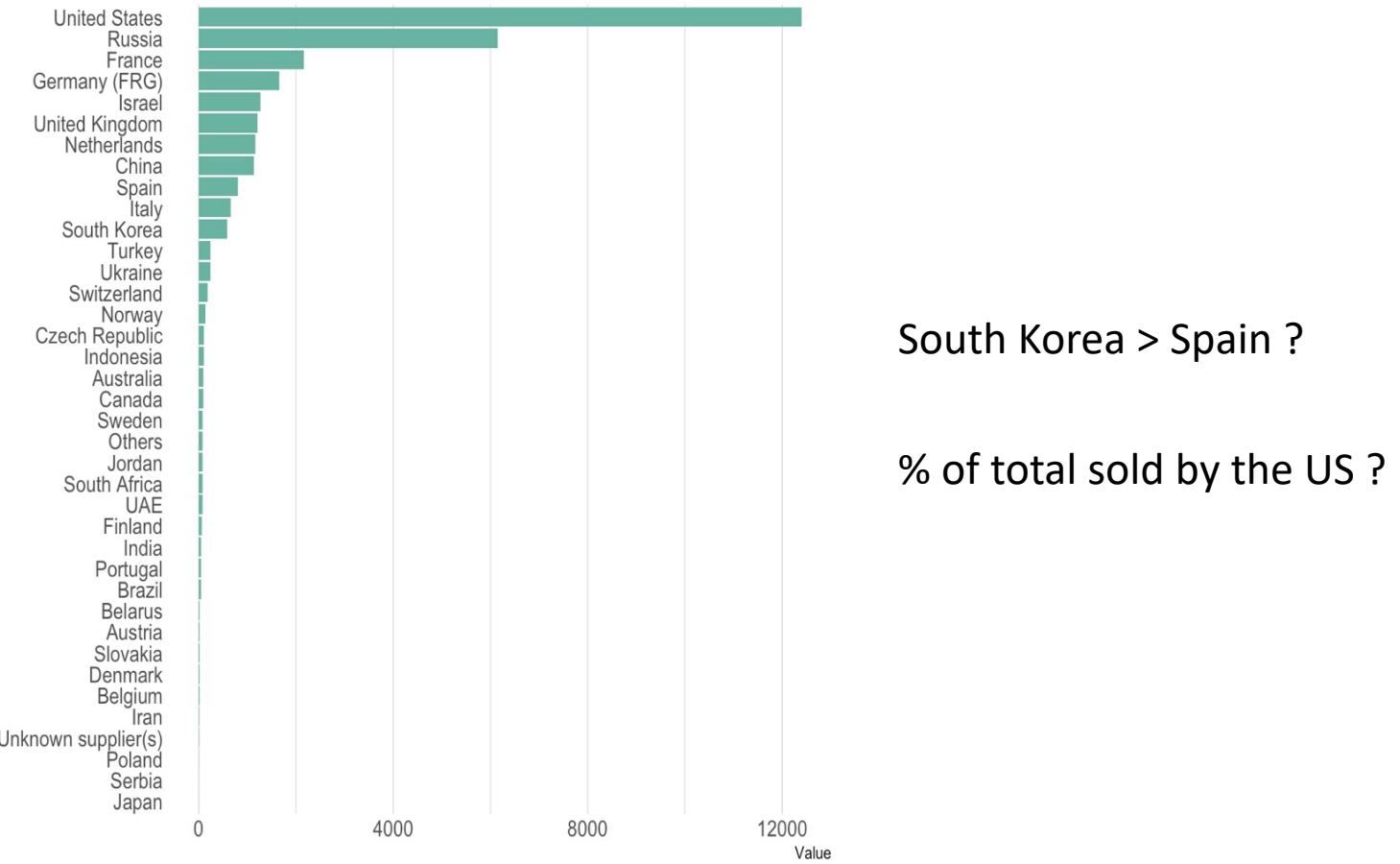


Who sells more weapons ?

Country	Value
United States	12394
Russia	6148
Germany (FRG)	1653
France	2162
United Kingdom	1214
China	1131

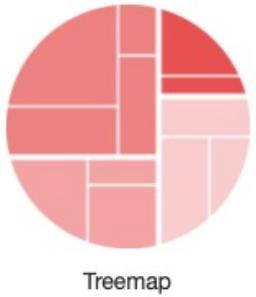


Barplot



Who sells more weapons ?

Country	Value
United States	12394
Russia	6148
Germany (FRG)	1653
France	2162
United Kingdom	1214
China	1131



United States

Russia

Israel

United Kingdom

Netherlands

China

Spain

South Korea

Turkey

Ukraine

France

Germany (FRG)

Italy

Norway

Switzerland

Indonesia

Australia

Canada

UAE

South Africa

Sweden

Jordan

Pakistan

Portugal

Brazil

Czech Republic

Others

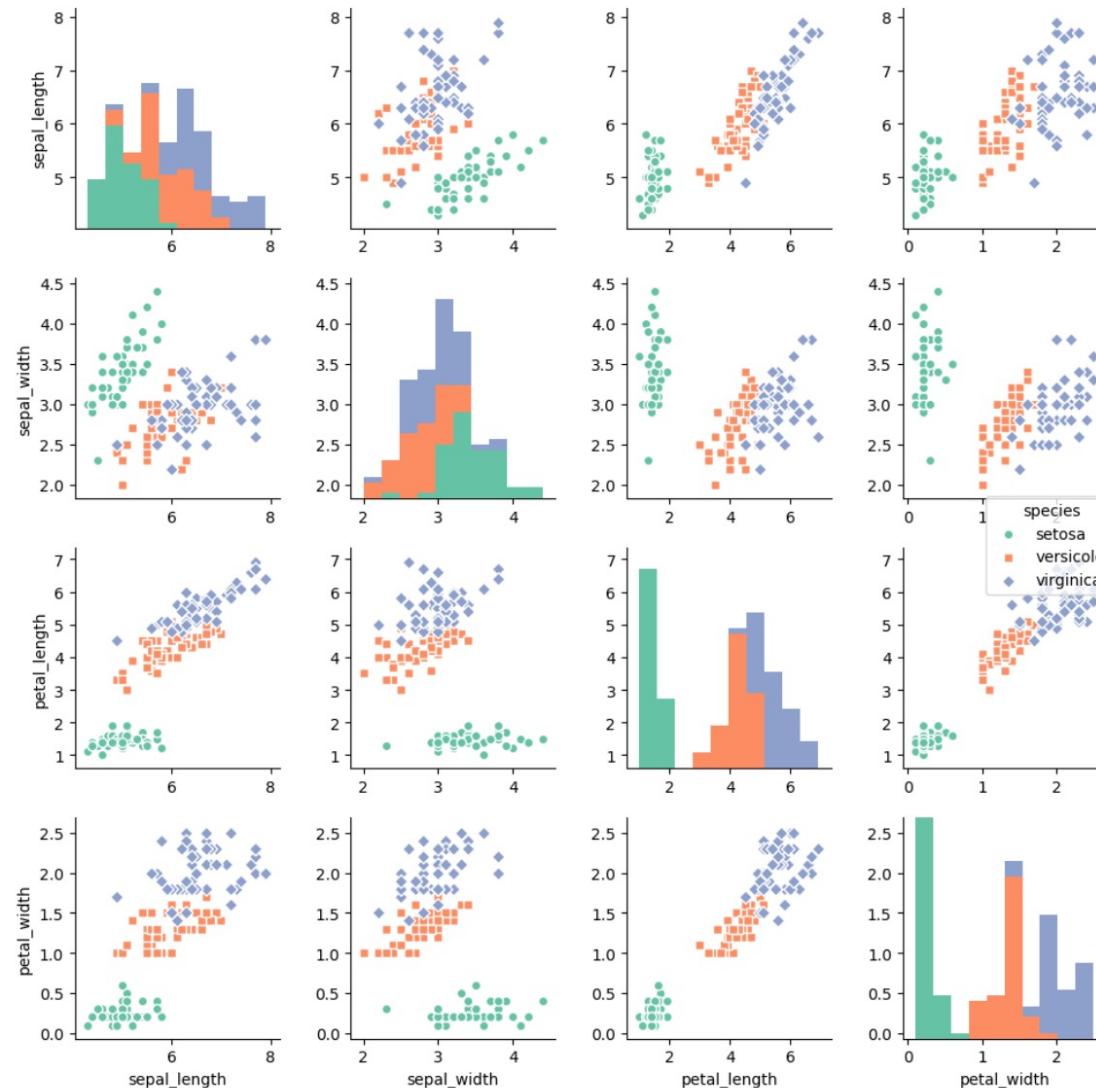
India

Others

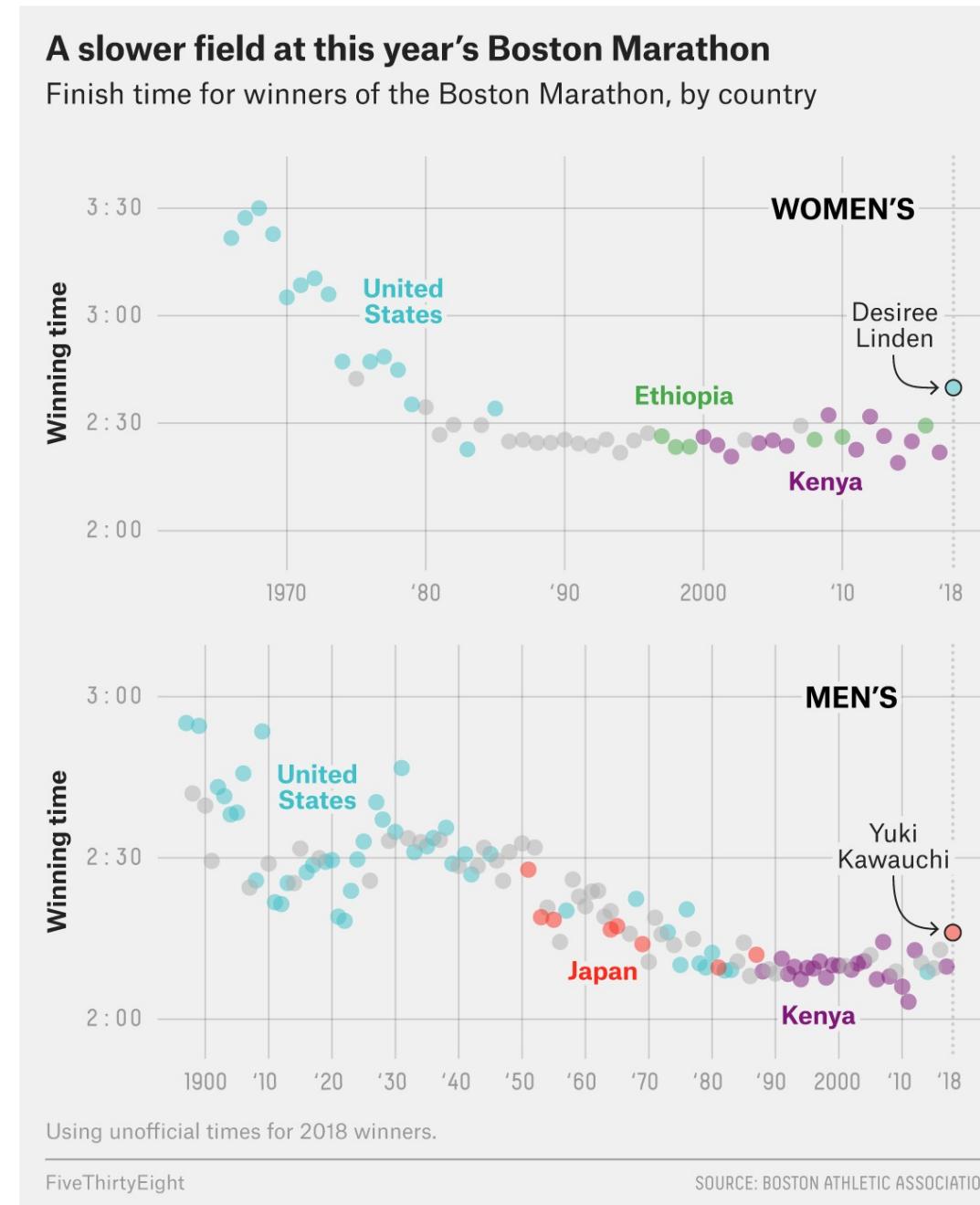
South Korea > Spain ?

% of total sold by the US ?

Question – **Exploratory Data Analysis** – Reader – Data – Interactivity – Eye catching



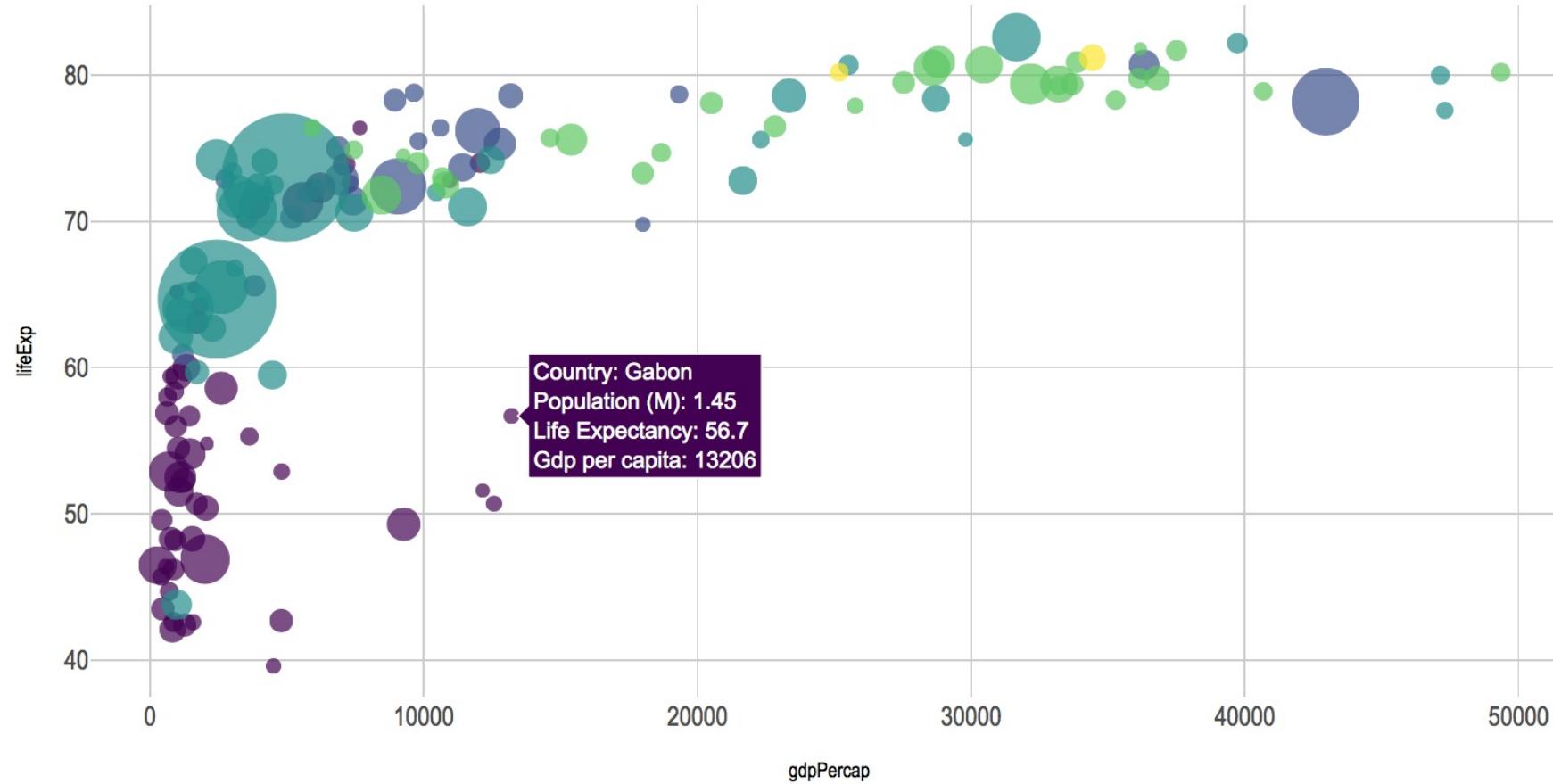
- Communicate result
- Show something specific
- Tells a story



Question – Explo/Expla – **Reader** – Data – Interactivity – Eye catching



Question – Explo/Expla – Reader – Data – **Interactivity** – Eye catching

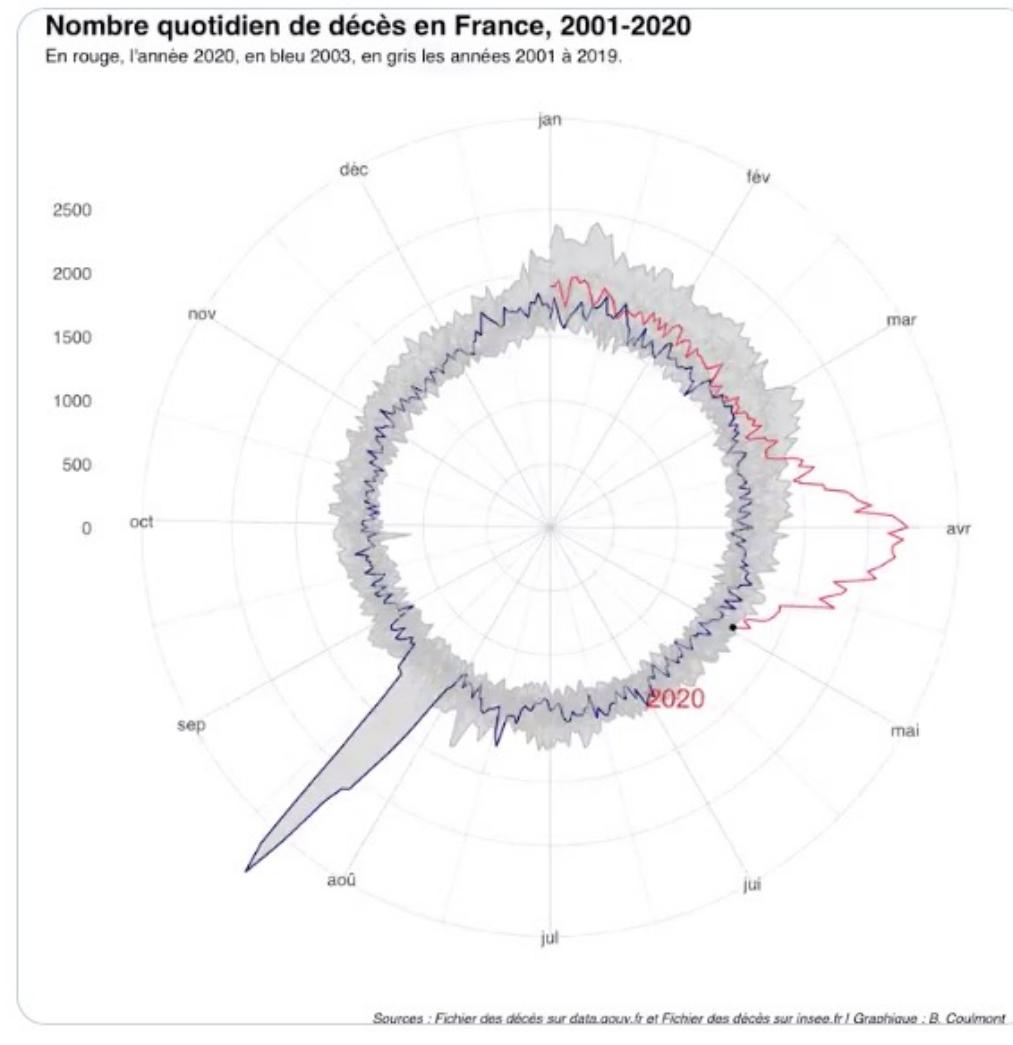


Question – Explo/Expla – Reader – Data – Interactivity – Eye catching





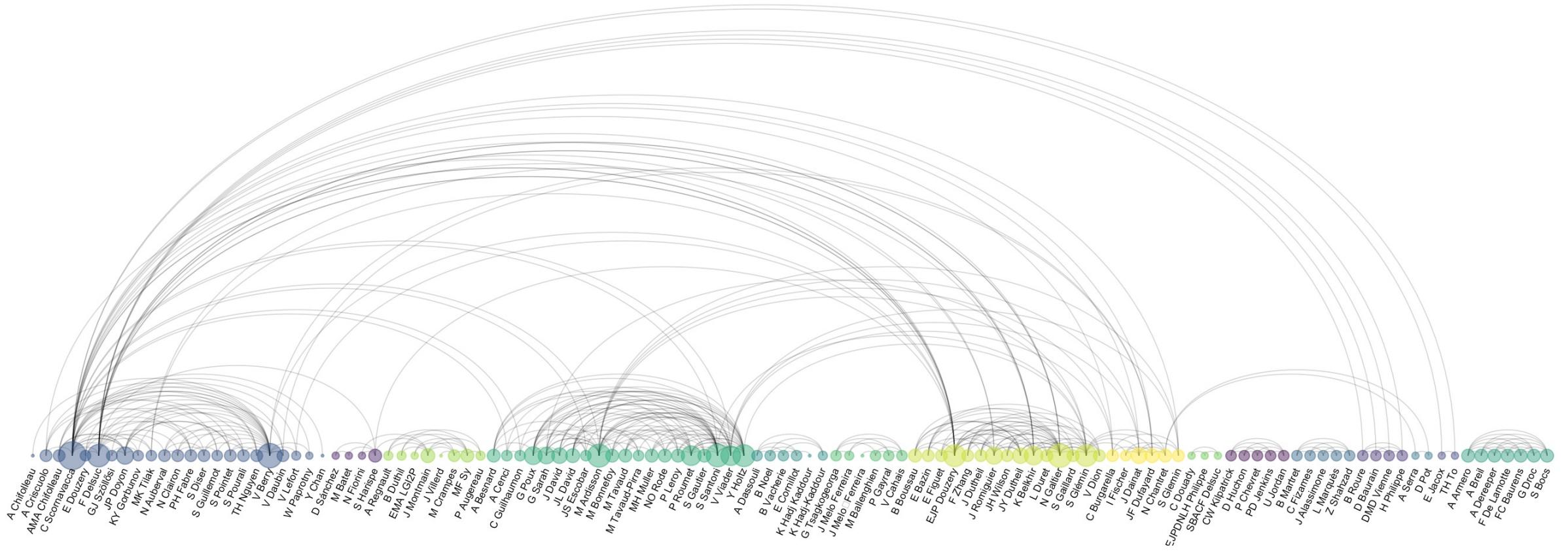
Nombre quotidien de décès en France, 2001-2020, en version animée, coordonnées polaires.



[Link](#)

12:43 PM · 2 déc. 2020 · Twitter Web App

6,5 k Retweets **1 k** Tweets cités **16,1 k** J'aime



Co-authorship network of a researcher



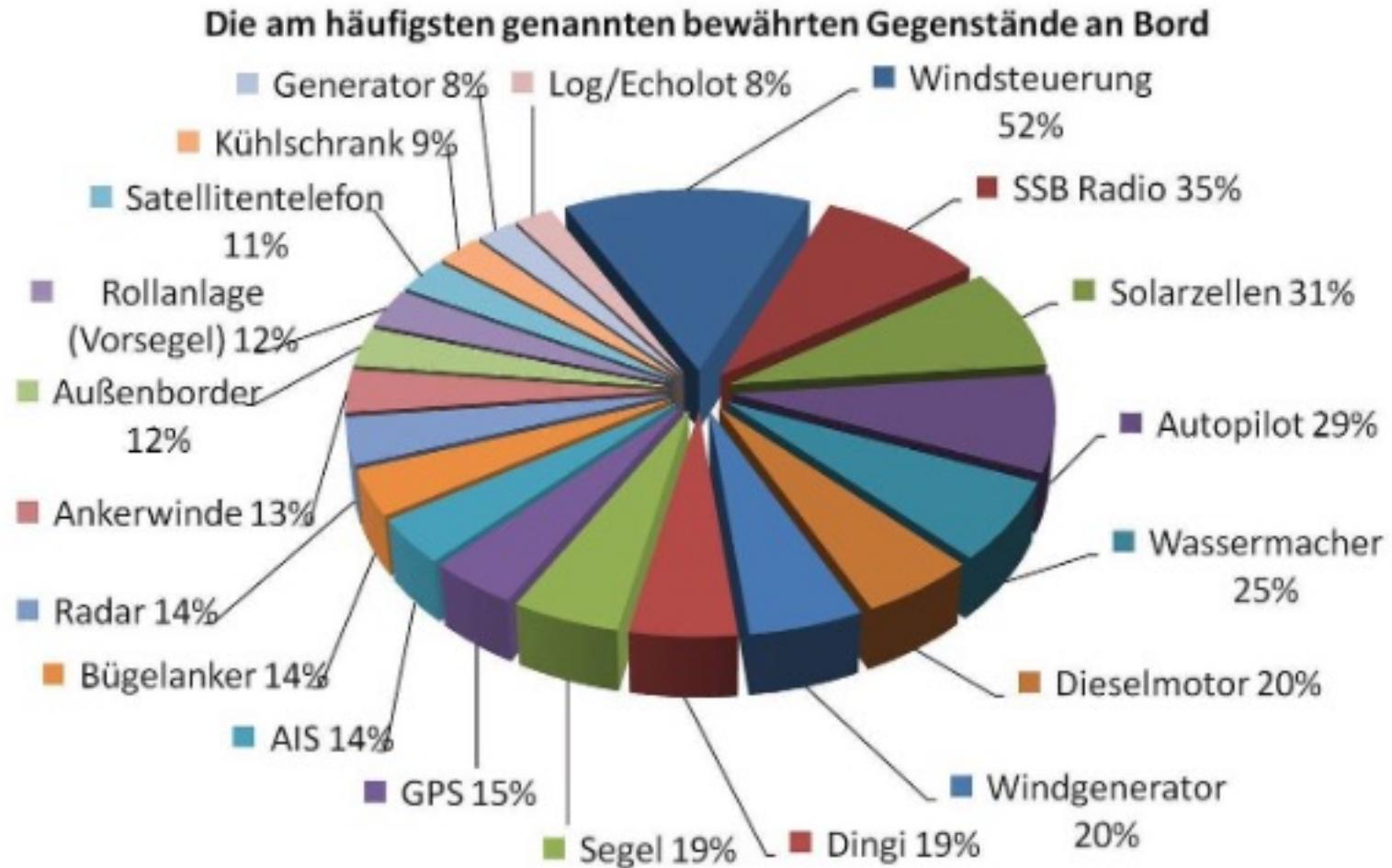
Where surfers travel.

data-to-viz.com | NASA.gov | 10,000 #surf tweets recovered

WHAT YOU SHOULD **NOT** DO

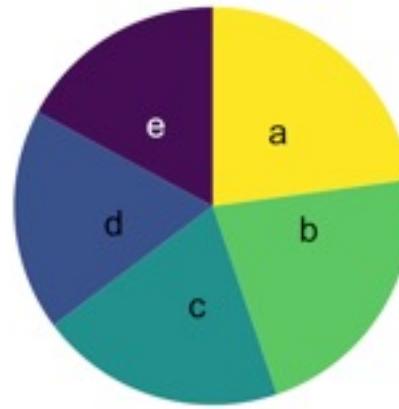
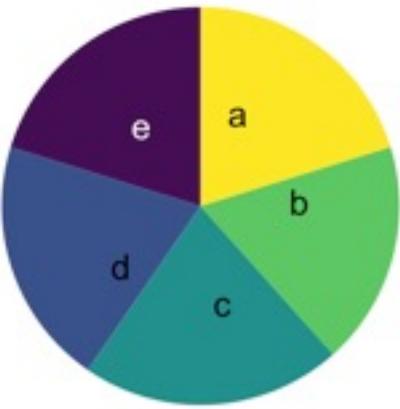
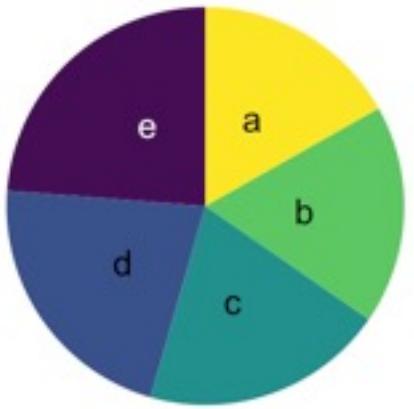
A gallery of most common caveats

What's wrong with
this chart?



Source: [WTF Visualizations](#)

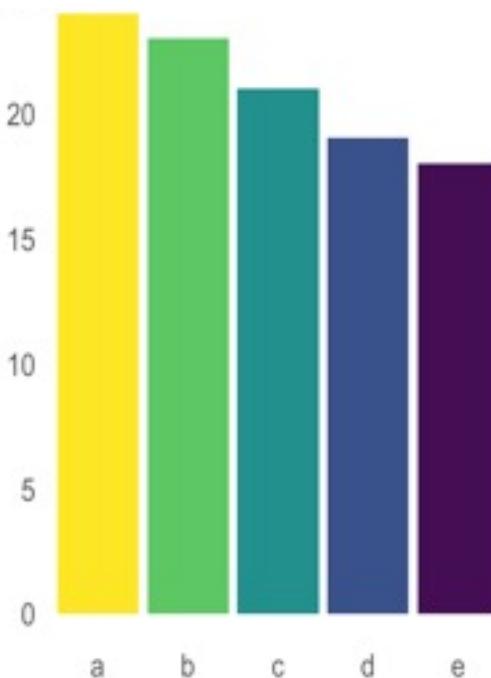
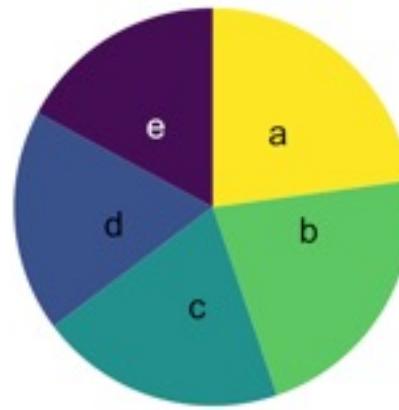
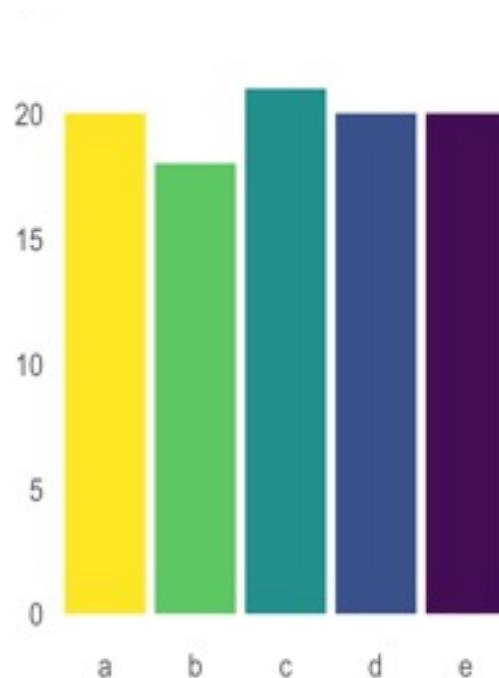
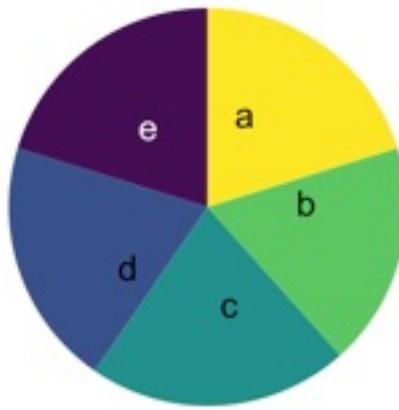
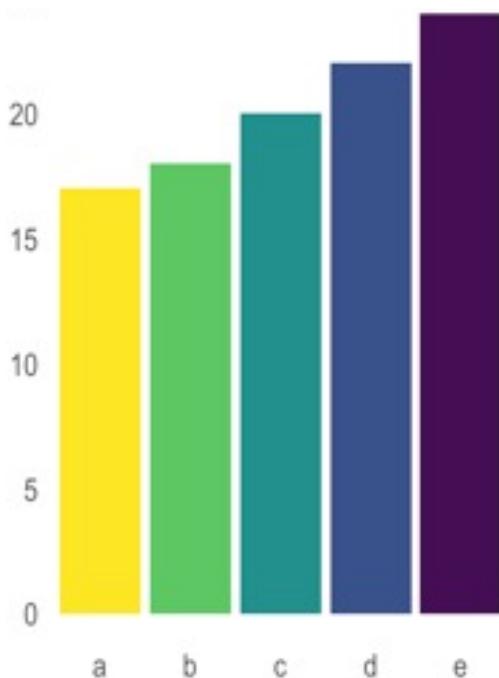
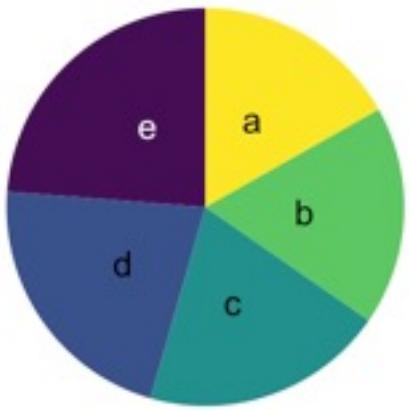
What's wrong
with pie chart?



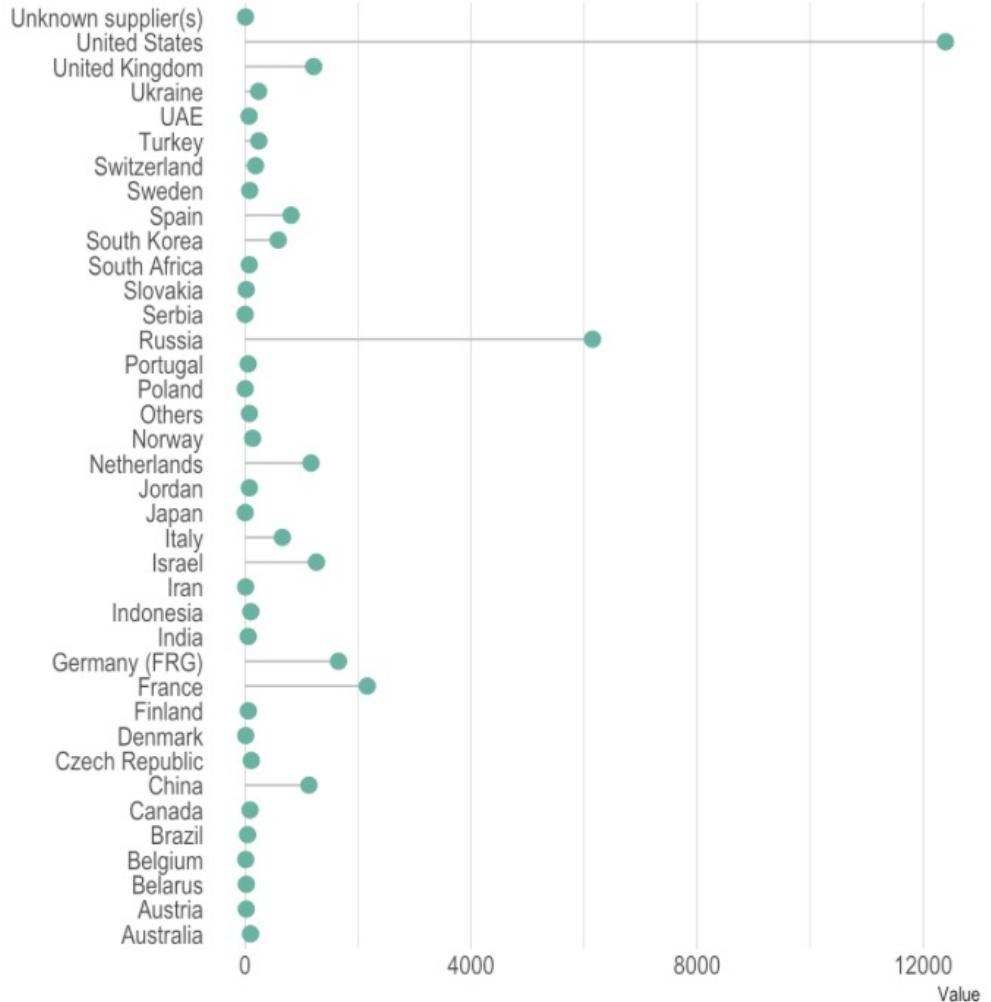
What can you see?

What's wrong with
pie chart?

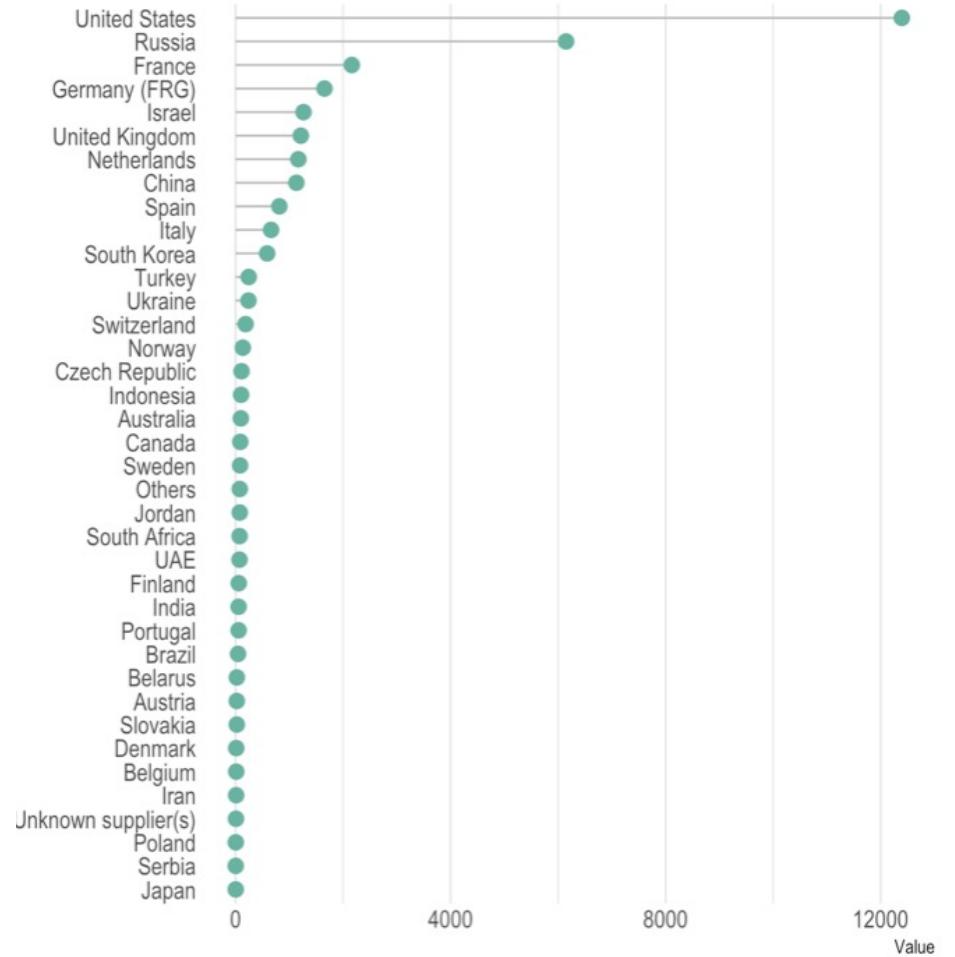
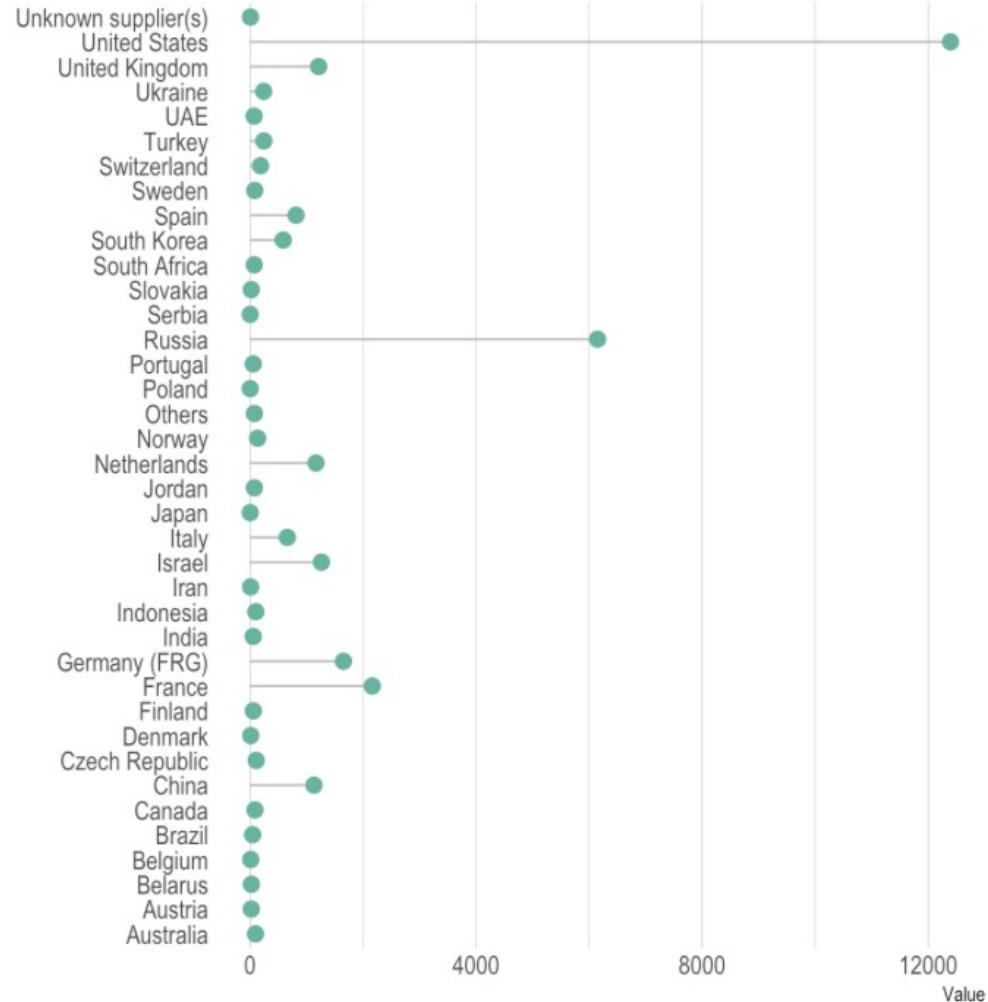
It is hard to
distinguish angles



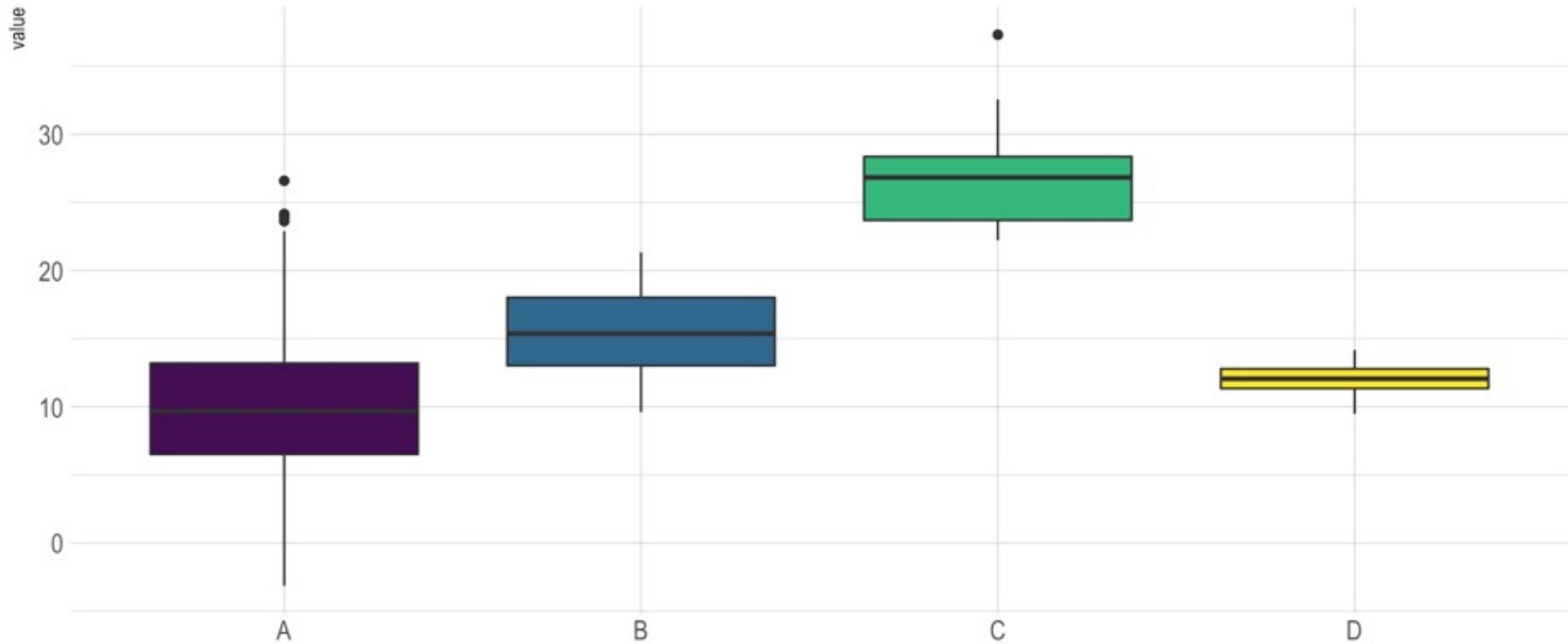
What could be better
here?



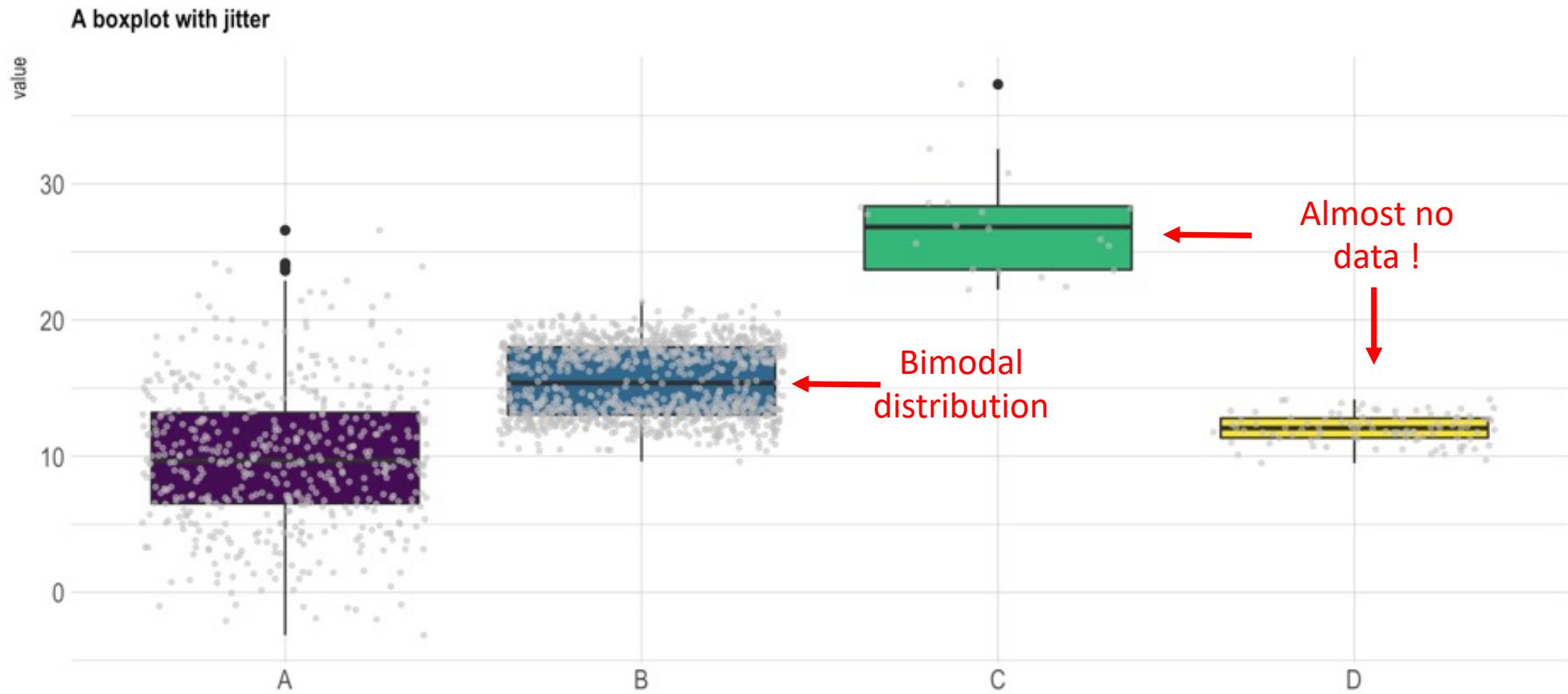
Order your data



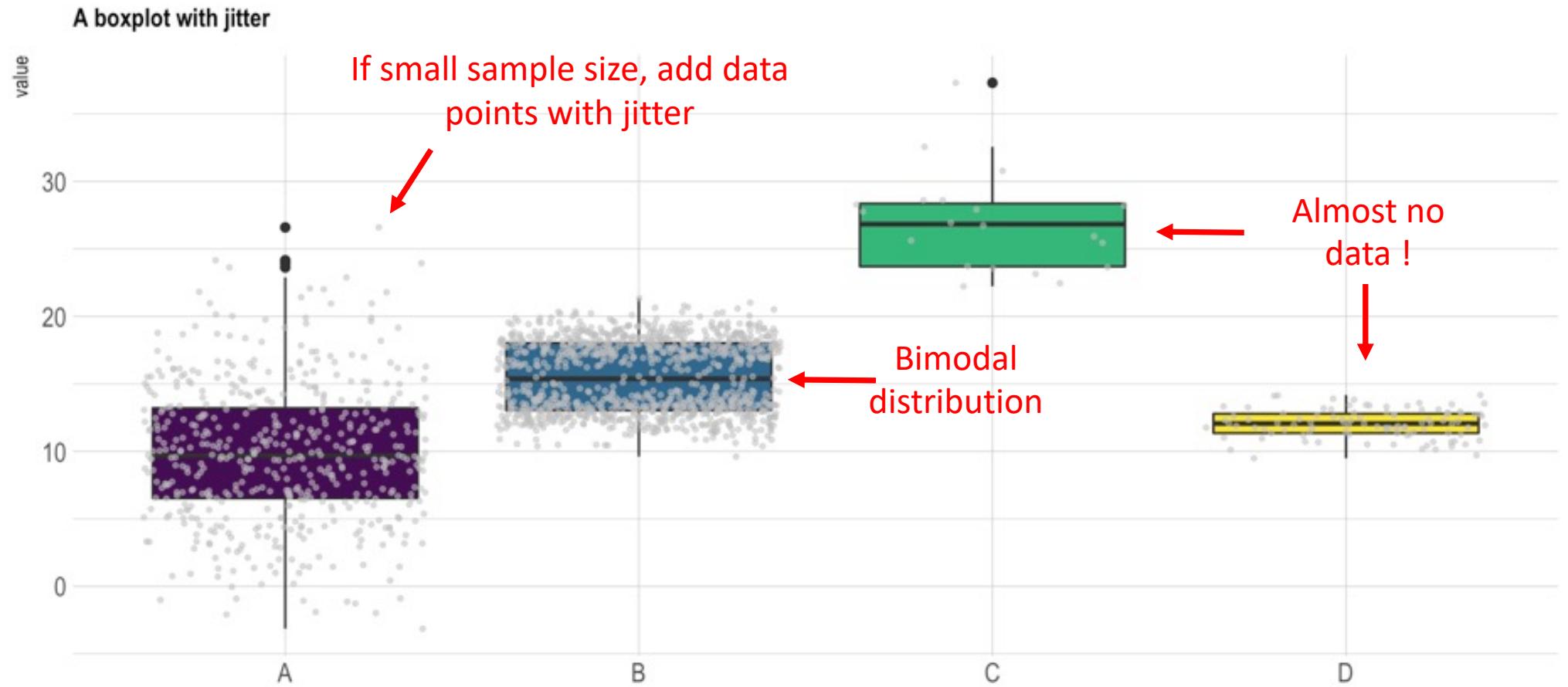
Anything wrong here?



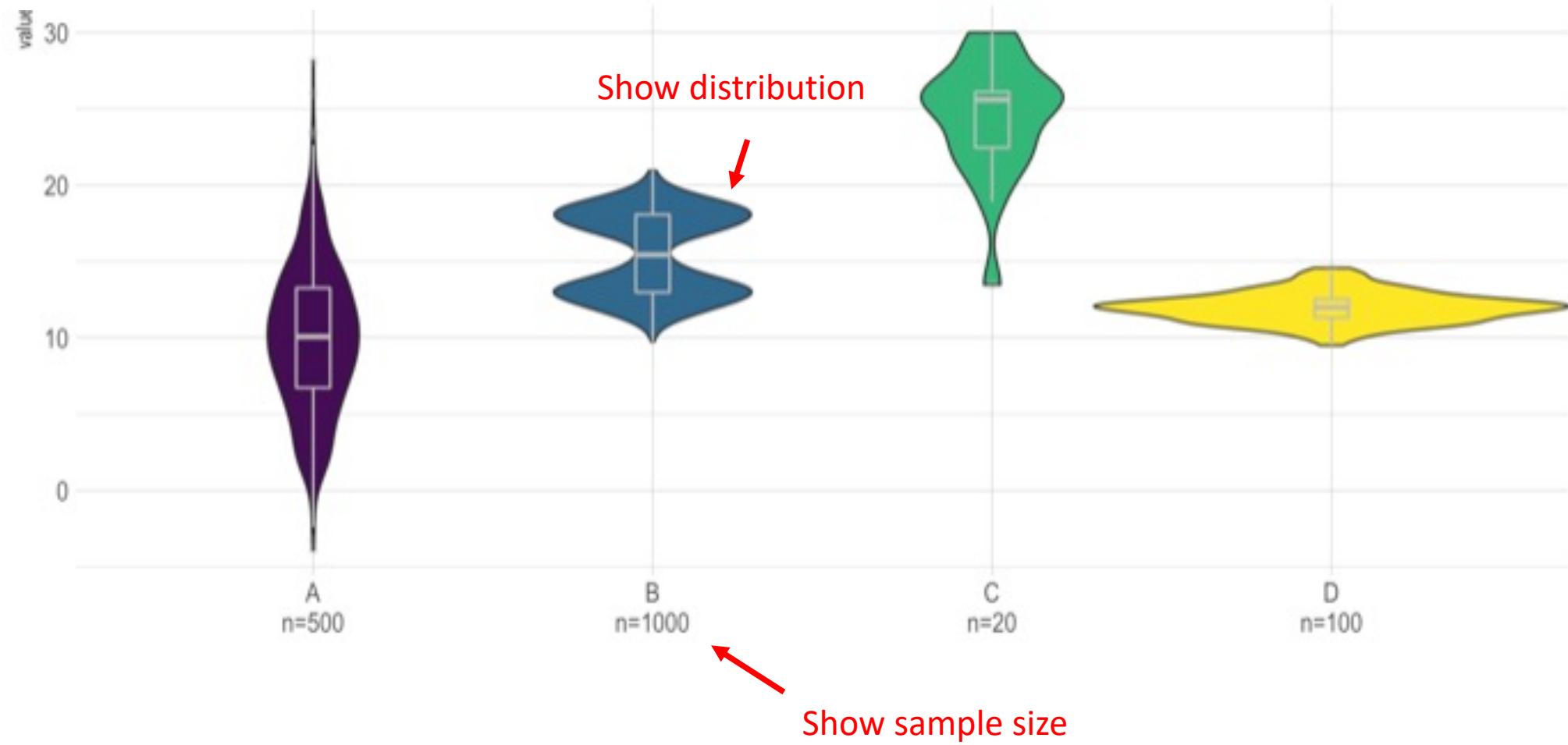
Boxplot = hide information



Boxplot = hide information



If big sample size, use violin plot





AMERICA'S ECONOMY

2010 GROSS DOMESTIC PRODUCT



United States
\$14.6 TRILLION



China \$5.7 TRILLION



Japan \$5.3 TRILLION

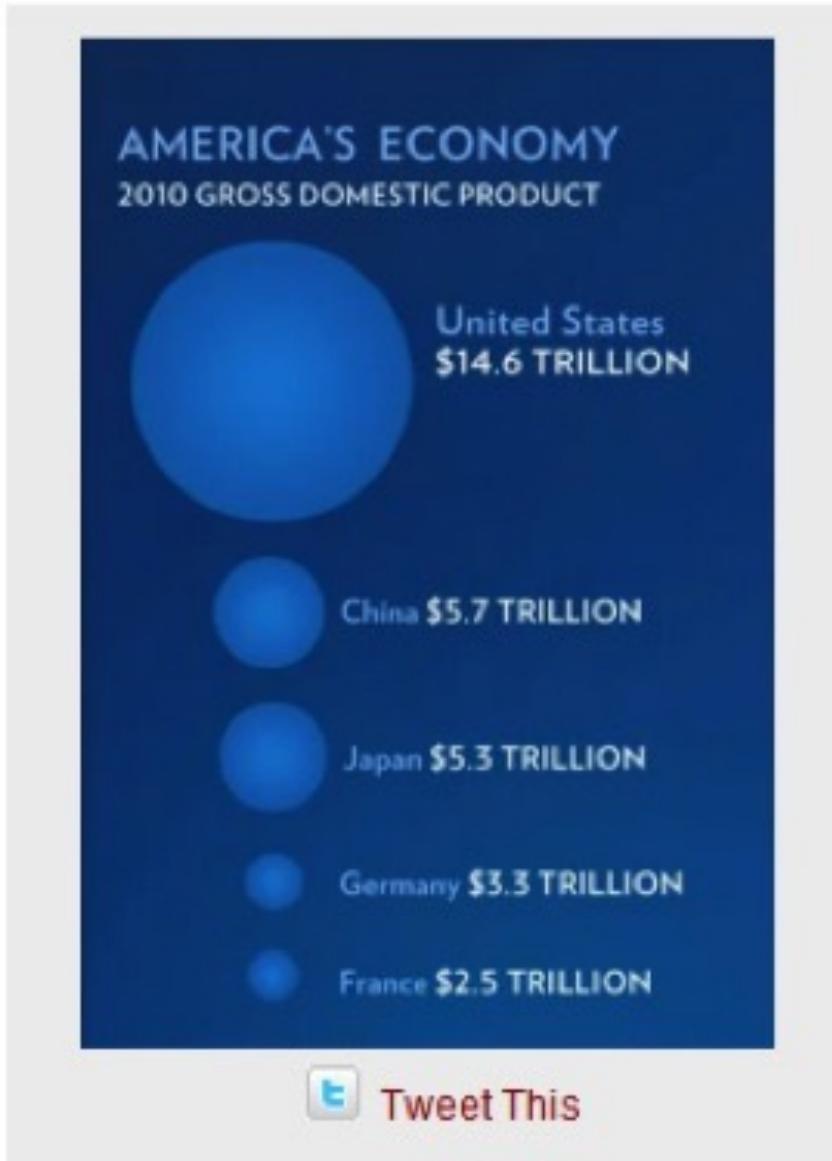


Germany \$3.3 TRILLION

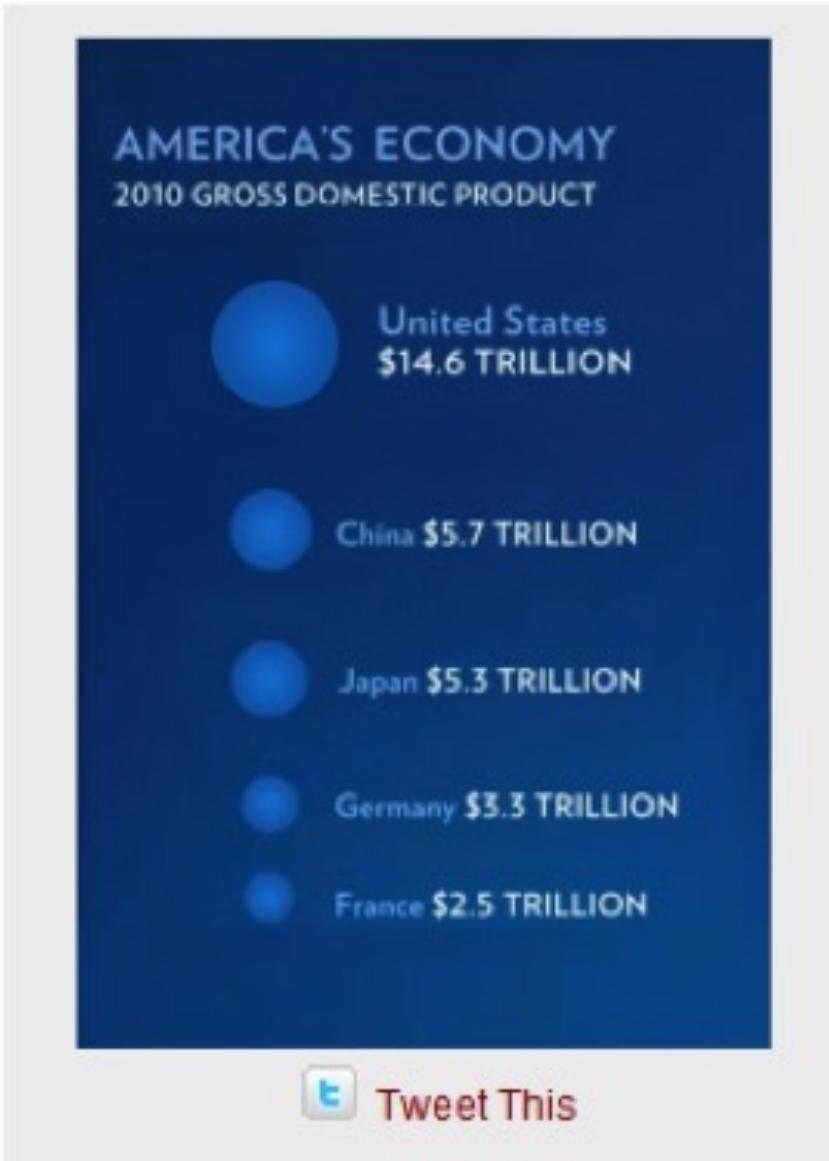


France \$2.5 TRILLION

Size = radius

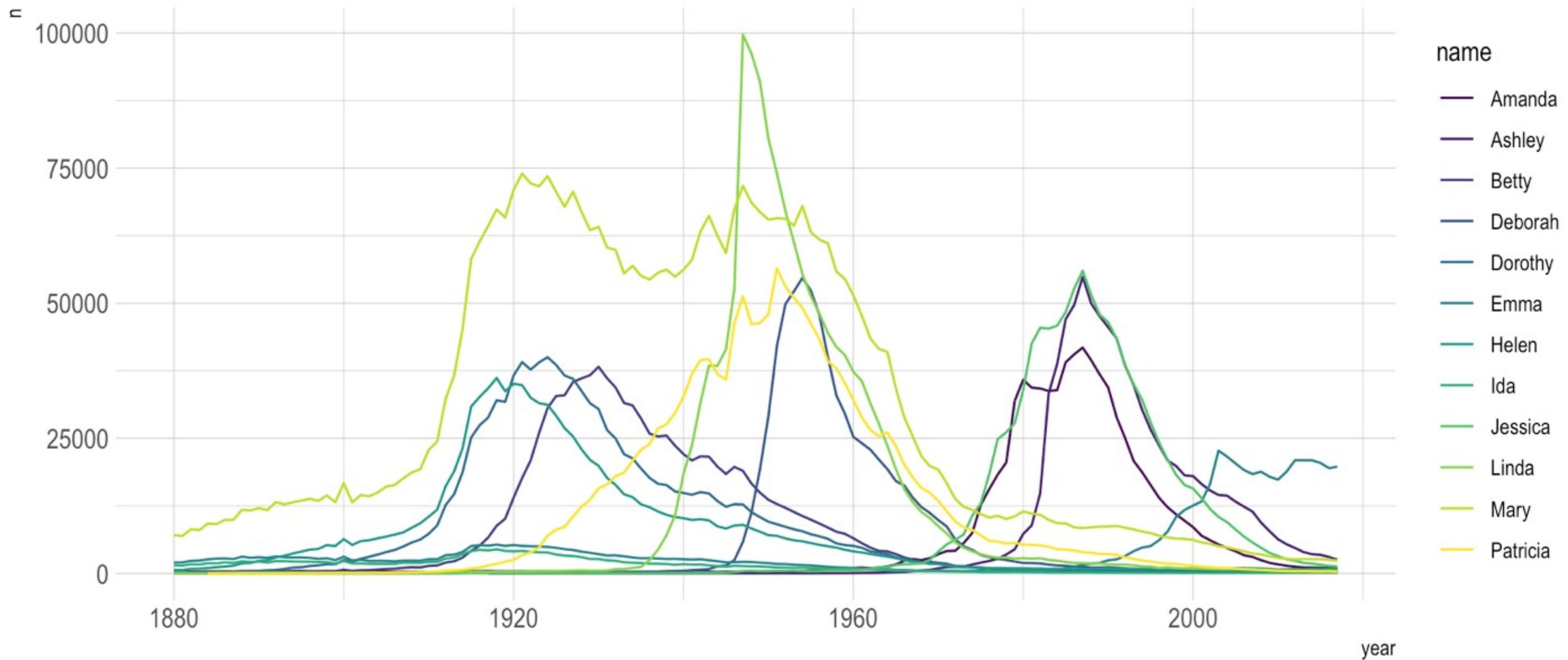


Size = area

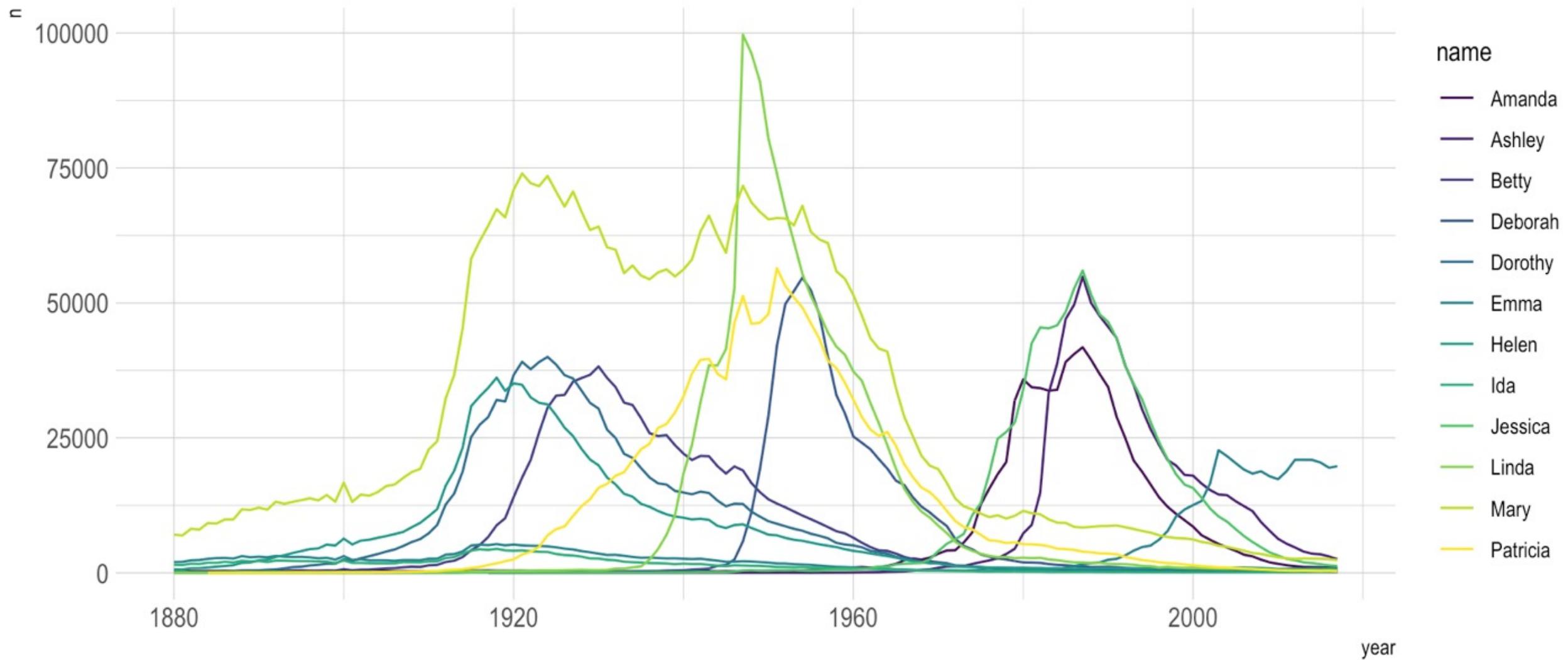


Source: [Fast Fedora blog](#)

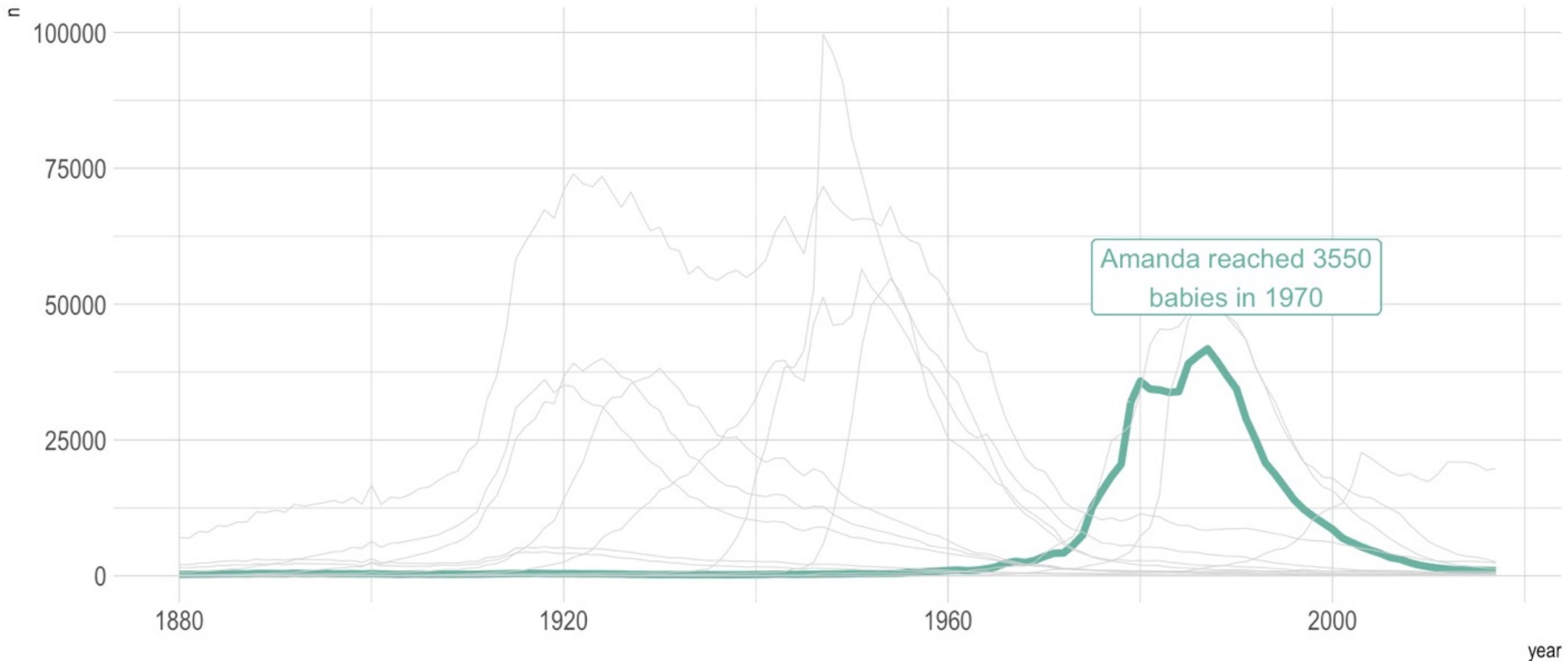
A spaghetti chart of baby names popularity



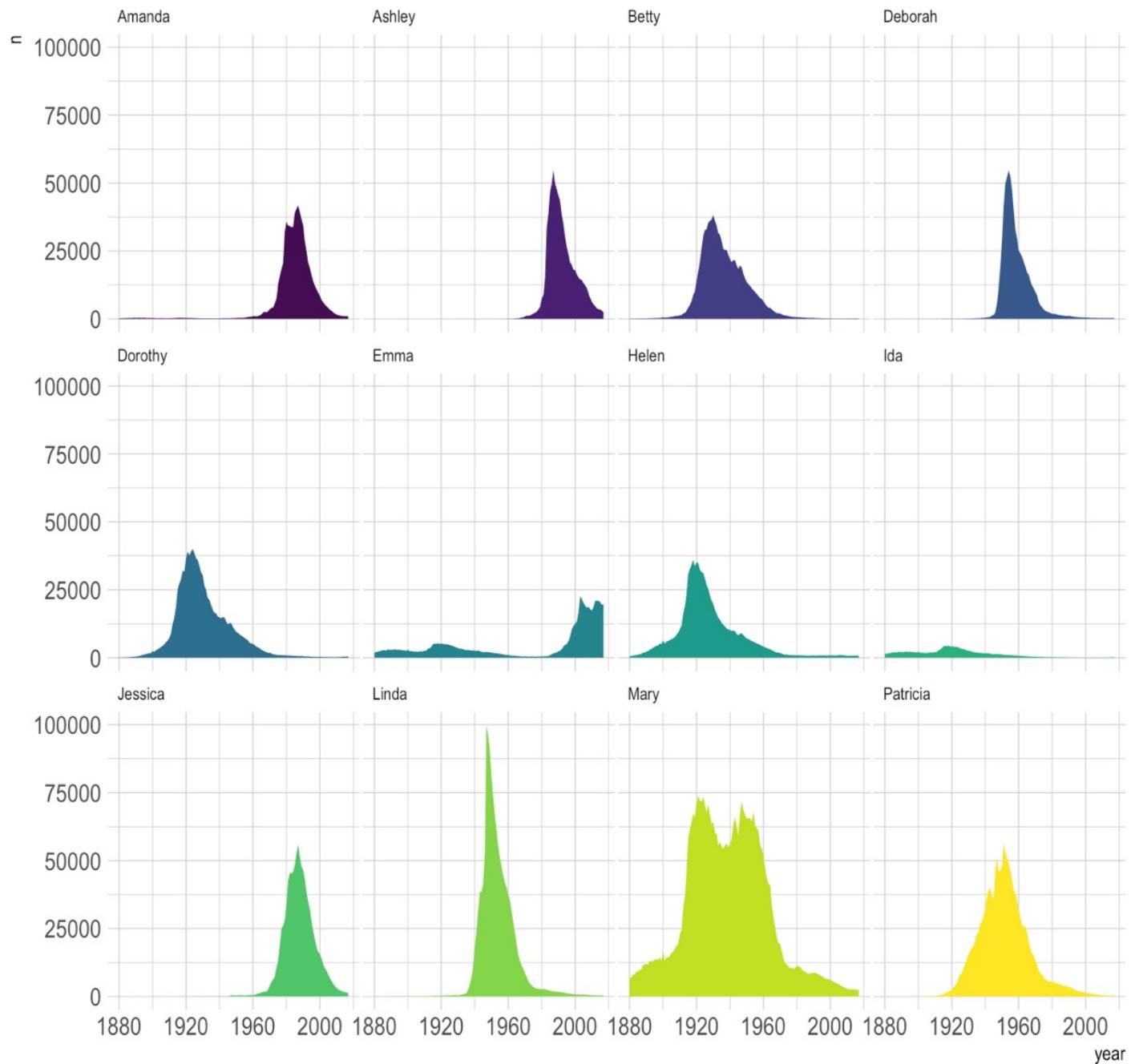
A spaghetti chart of baby names popularity



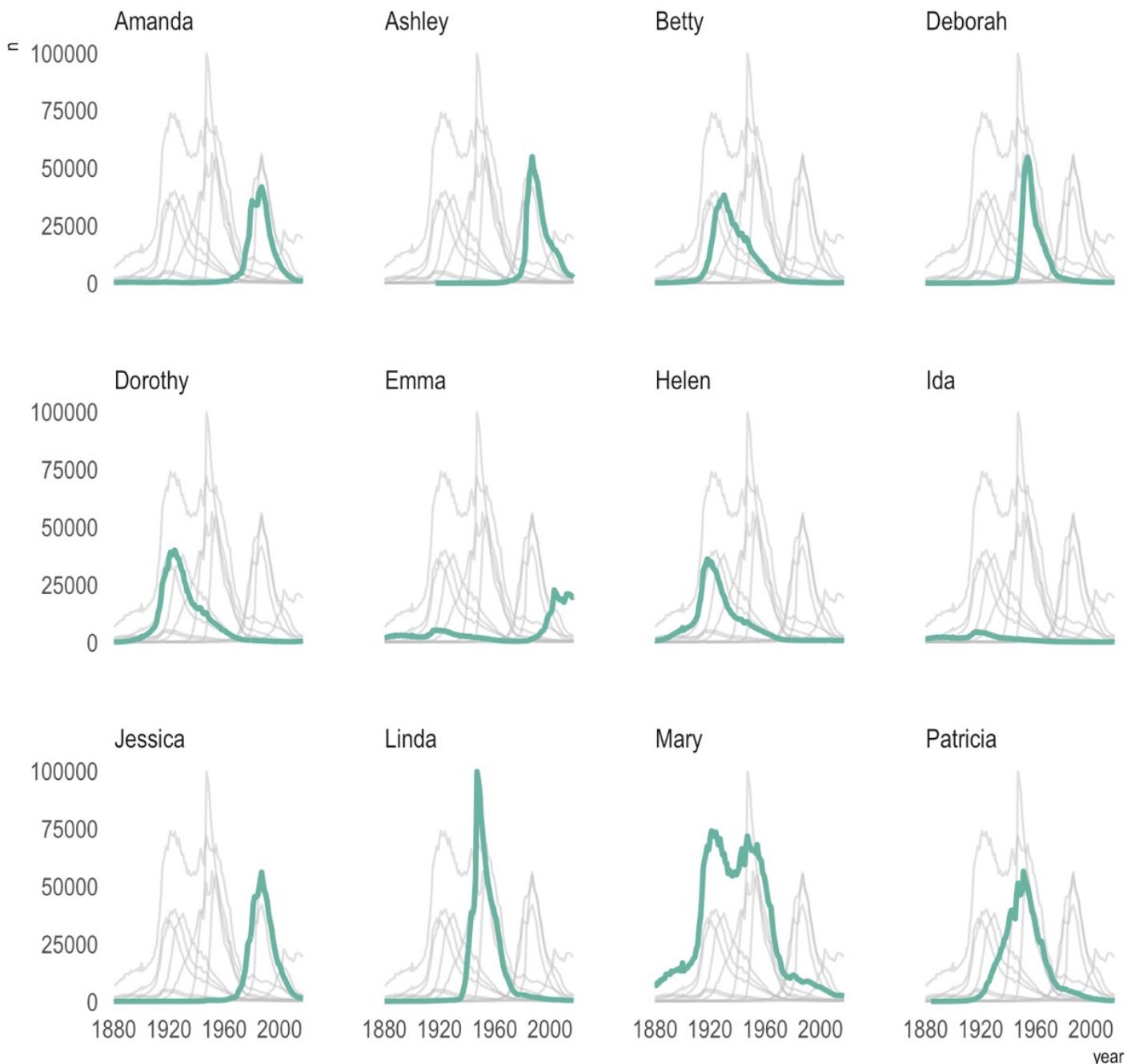
Popularity of American names in the previous 30 years

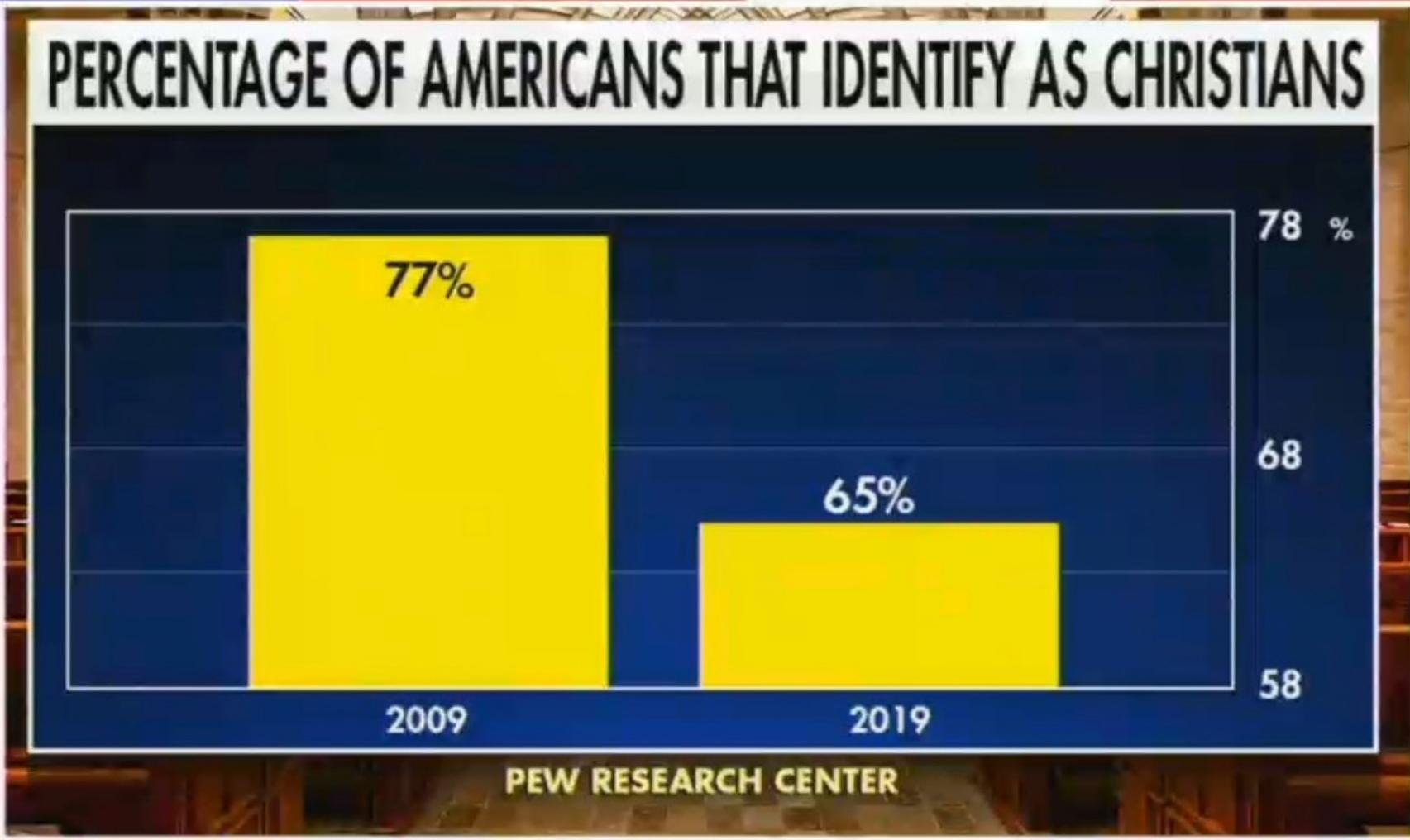


Popularity of American names in the previous 30 years



A spaghetti chart of baby names popularity





WE USED TO BE AN ENTHUSIASTICALLY CHRISTIAN NATION
TUCKER CARLSON • TONIGHT •

Data to Viz

A collection of
dataviz caveats

[Data-to-viz.com/caveats](https://www.dat-to-viz.com/caveats)



Order your data

When displaying the value of several entities, ordering them makes the graph much more insightful.



To cut or not to cut?

Cutting the Y-axis is one of the most controversial practice in data viz. See why.



The spaghetti chart

A line graph with too many lines becomes unreadable: it is called a spaghetti graph.



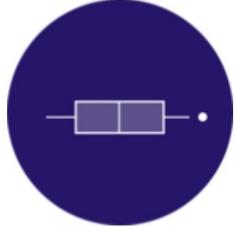
Pie chart

The human eye is bad at reading angles. See how to replace the most criticized chart ever.



Play with histogram bin size

Always try different bin sizes when you build a histogram, it can lead to different insights.



Do boxplots hide information?

Boxplots are a great way to summarize a distribution but hide the sample size and their distribution.



The problem with error bars

Barplots with error bars must be used with great care. See why and how to replace them.



Too many distributions.

If you need to compare the distributions of many variables, don't clutter your graphic.

HOW TO DO IT

The R and Python graph galleries

Design the idea



20 %

Build it



80 %



Easy

Hard

Limited

Flexible

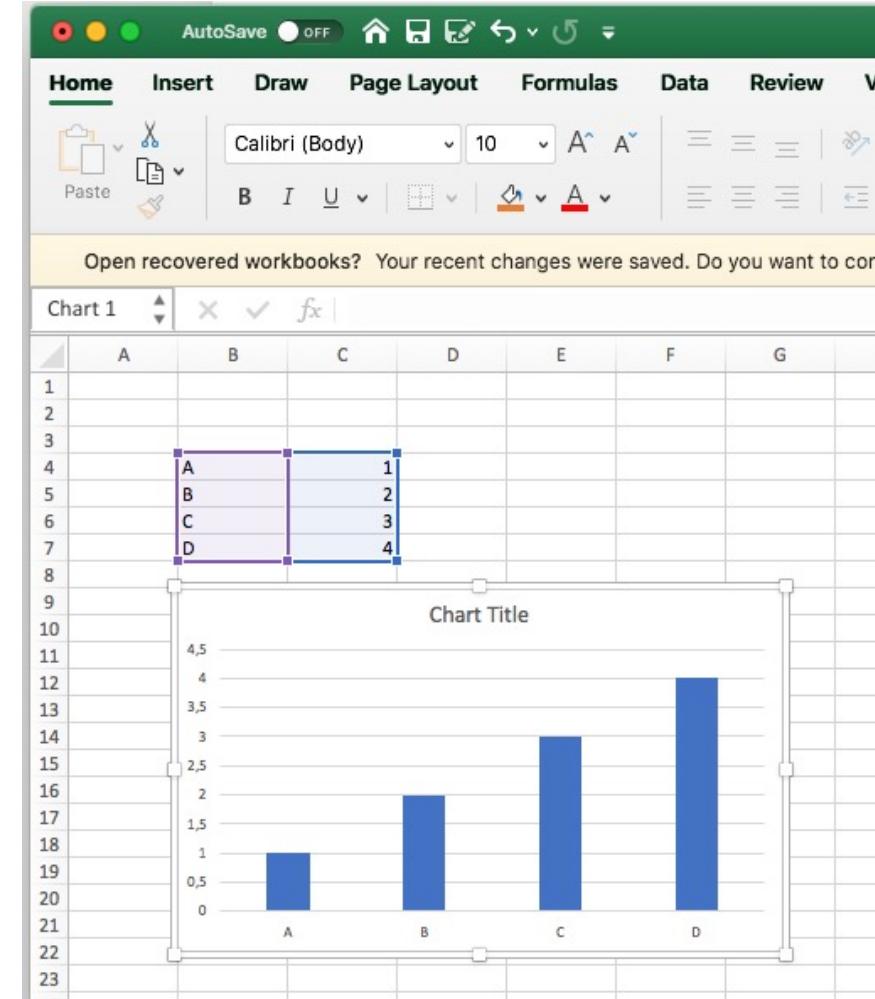




Easy

Excel

Limited

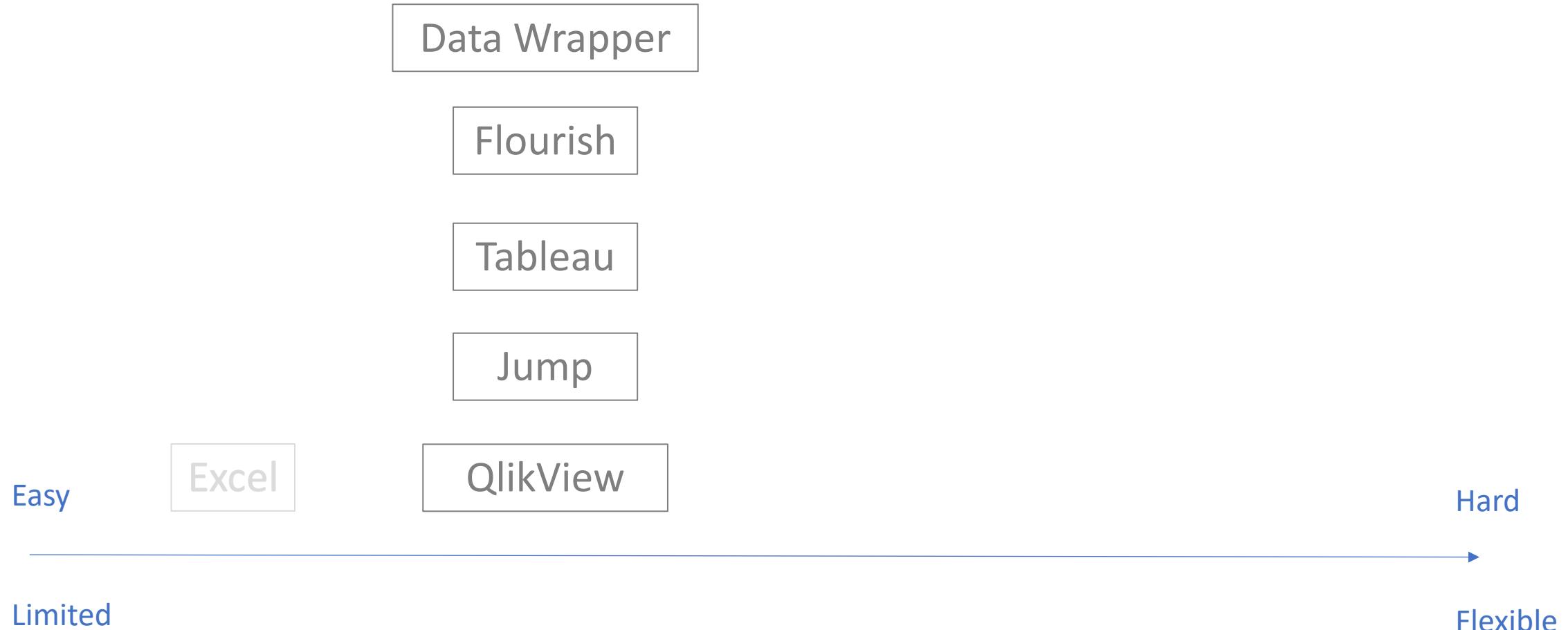


Hard

Flexible

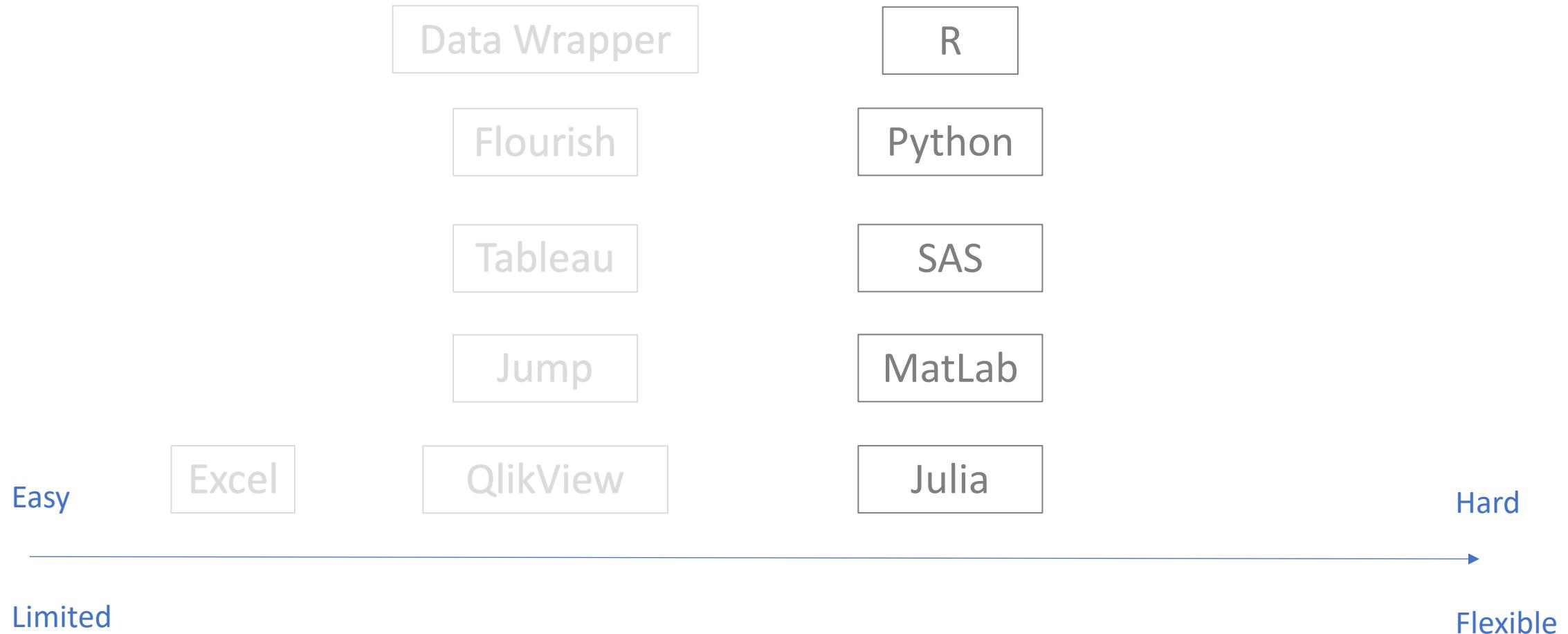


[Data Wrapper demo](#)





[R graph gallery demo](#)



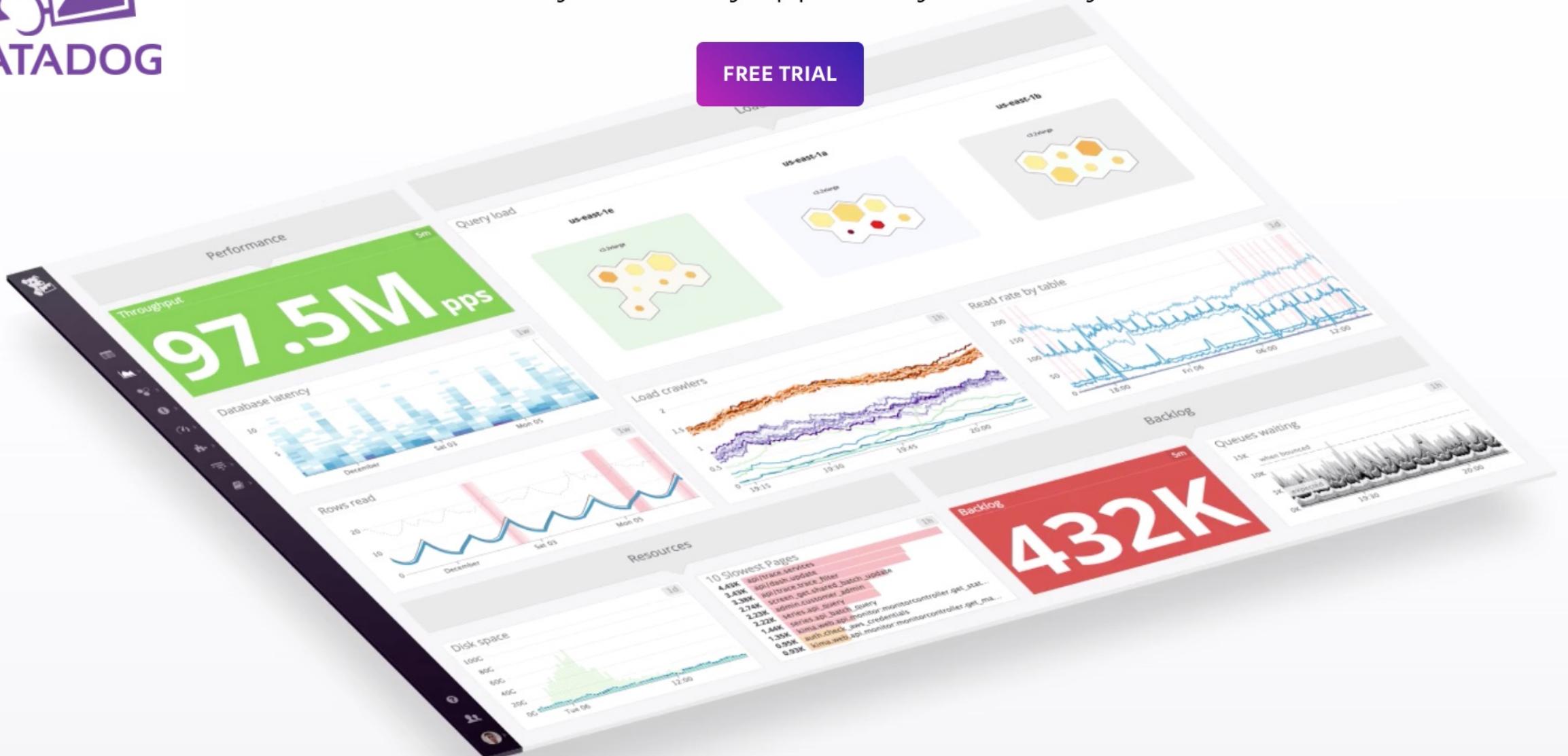




Modern monitoring & analytics

See inside any stack, any app, at any scale, anywhere.

FREE TRIAL





from Data to *Viz*



[**R-graph-gallery.com**](#)

[**Python-graph-gallery.com**](#)

[**D3-graph-gallery.com**](#)



KANTAR
Information is Beautiful
Awards

Data-to-viz.com



@R_Graph_Gallery



github.com/holtzy/Talk



Yan.holtz.data@gmail.com



www.yan-holtz.com