

Data types and standard visualizations

Anvar Kurmukov, 2022

Exploratory data analysis (EDA)

Exploratory Data Analysis (EDA) is an approach for data analysis without statistical model or formulated prior hypothesis

- Maximize insight into a data set
- Uncover underlying structure
- Detect missing data
- Detect outliers and anomalies
- Rank important factors
- Perform sanity check

Descriptive statistics: computing simple summary statistics such as mean, median, standard deviation, plotting box plots, histograms

Visualization: plotting the raw data - data traces, scatter plots, frequency plots, probability plots, multivariate plots

Data types

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult_male	deck	embark_town	alive	alone
0	0	3	male	22.0	1	0	7.2500	S	Third	man	True	NaN	Southampton	no	False
1	1	1	female	38.0	1	0	71.2833	C	First	woman	False	C	Cherbourg	yes	False
2	1	3	female	26.0	0	0	7.9250	S	Third	woman	False	NaN	Southampton	yes	True
3	1	1	female	35.0	1	0	53.1000	S	First	woman	False	C	Southampton	yes	False
4	0	3	male	35.0	0	0	8.0500	S	Third	man	True	NaN	Southampton	no	True
5	0	3	male	NaN	0	0	8.4583	Q	Third	man	True	NaN	Queenstown	no	True
6	0	1	male	54.0	0	0	51.8625	S	First	man	True	E	Southampton	no	True
7	0	3	male	2.0	3	1	21.0750	S	Third	child	False	NaN	Southampton	no	False
8	1	3	female	27.0	0	2	11.1333	S	Third	woman	False	NaN	Southampton	yes	False

- TABLE ROWS = instances, examples, data points, observations, samples
- TABLE COLUMNS = attributes, features, variables
- Categorical data (= labels, nominal, ordinal [ordered], binary)
- Quantitative data (= numbers, discrete [integer], continues [real])

EDA tools Methods and approaches

Summary statistics

- min, max (range)
- mean, median (location)
- variance, standard deviation (dispersion)
- skewness (asymmetry)
- kurtosis (peakedness)

Visualization

- Bar plots
- Scatter plots
- Histograms
- Box plots
- Index chart

Mean

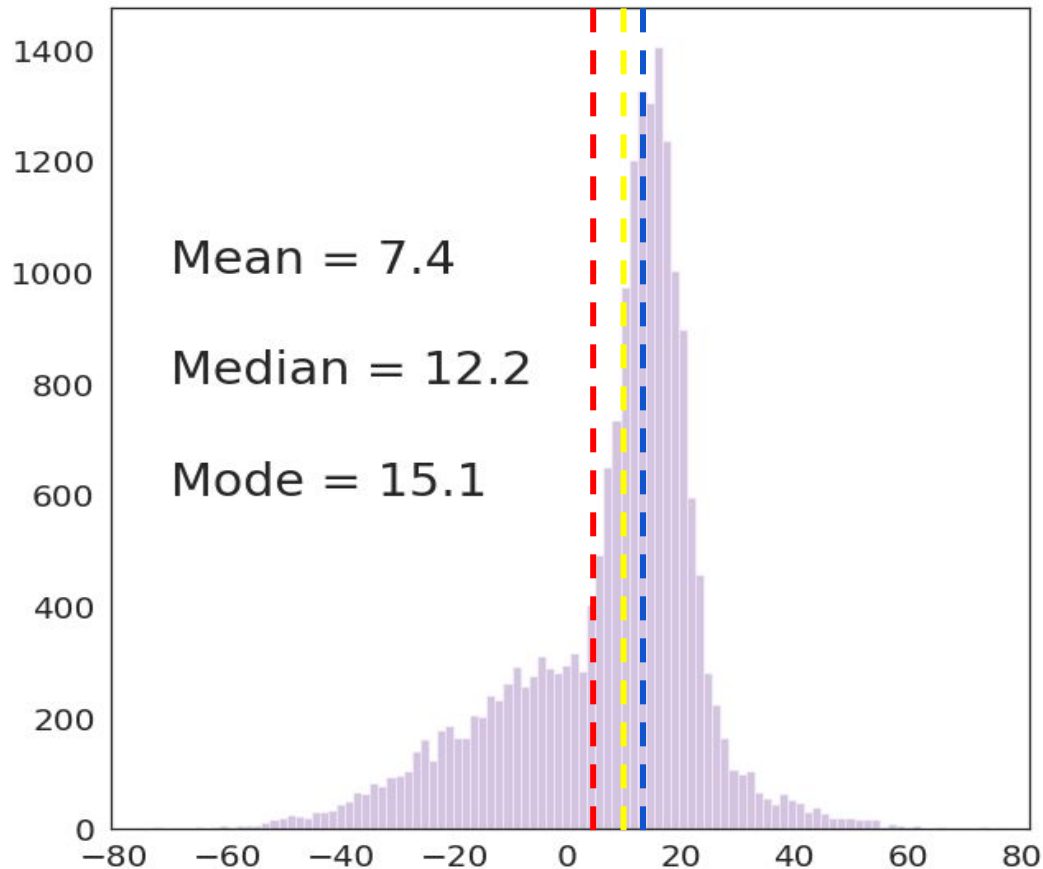
arithmetic mean

Median

50-th percentile

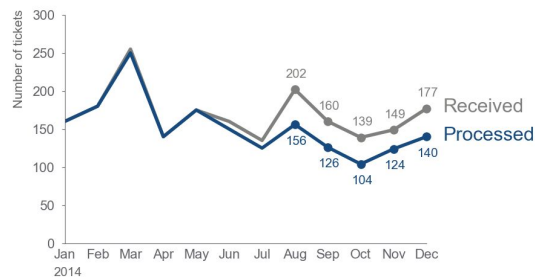
Mode

the most frequent value



Standard graphs

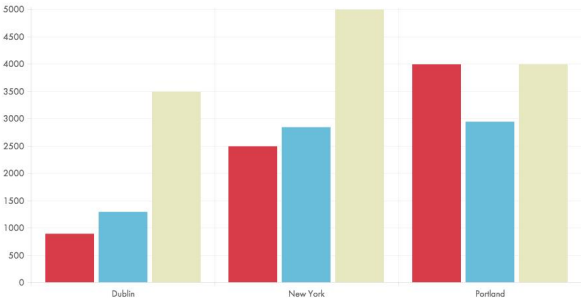
Ticket volume over time



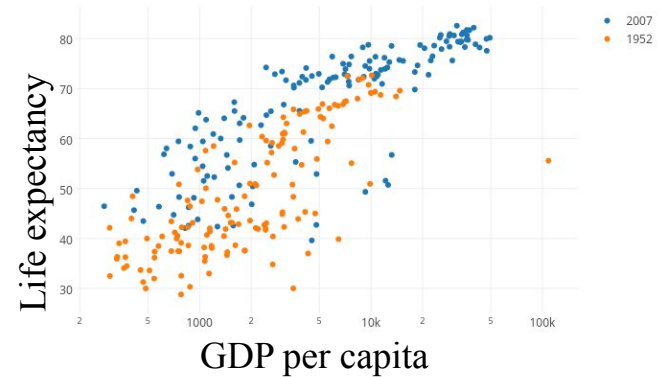
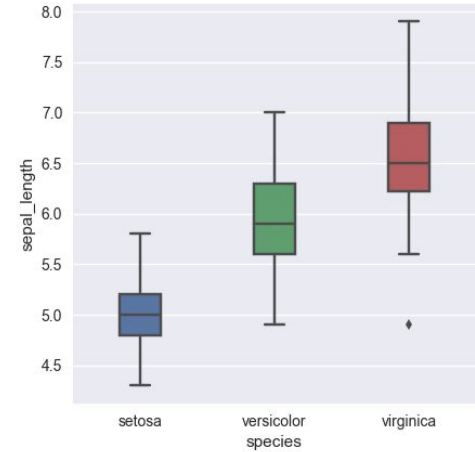
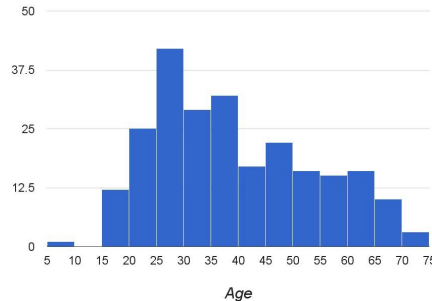
Data source: XYZ Dashboard, as of 12/31/2014

1. Index chart
2. Bar chart (pie chart)
3. Box plot
4. Scatter plot
5. Histogram

TICKET SALES BY LOCATION (JANUARY TO MARCH)



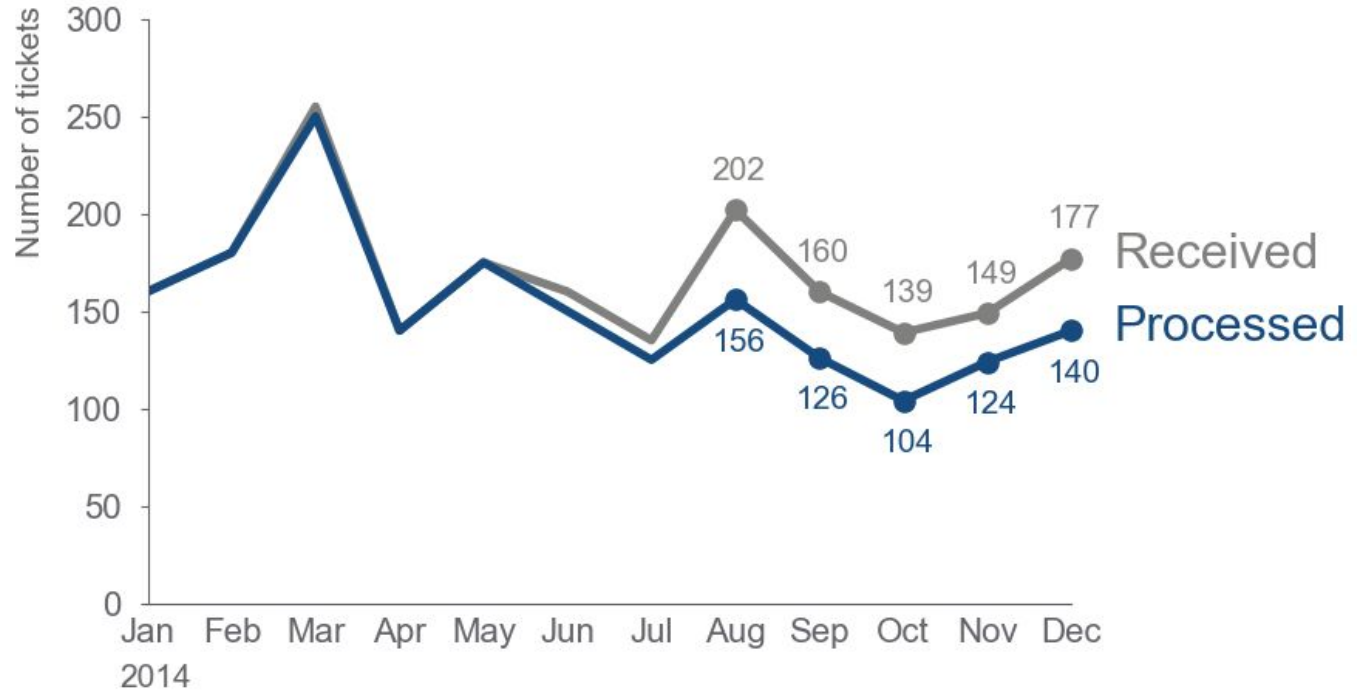
Respondents' age distribution (n=240)



Index chart

Index chart

Ticket volume over time



Data source: XYZ Dashboard, as of 12/31/2014

Index chart

Ticket volume over time



Data source: XYZ Dashboard, as of 12/31/2014

Time

Index chart

Ticket volume over time



Data source: XYZ Dashboard, as of 12/31/2014

Amount

Time

Histogram



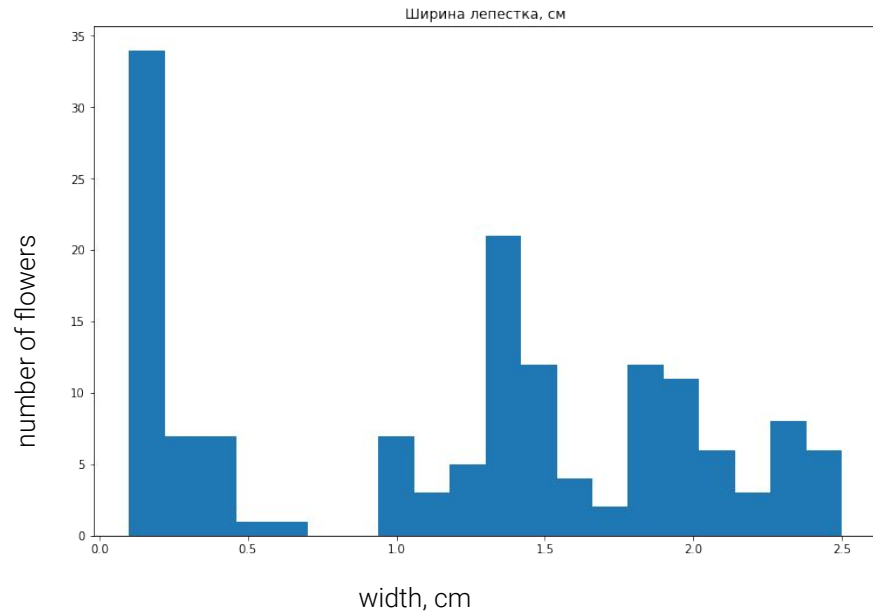
Ирис щетинистый
(*Iris setosa*)



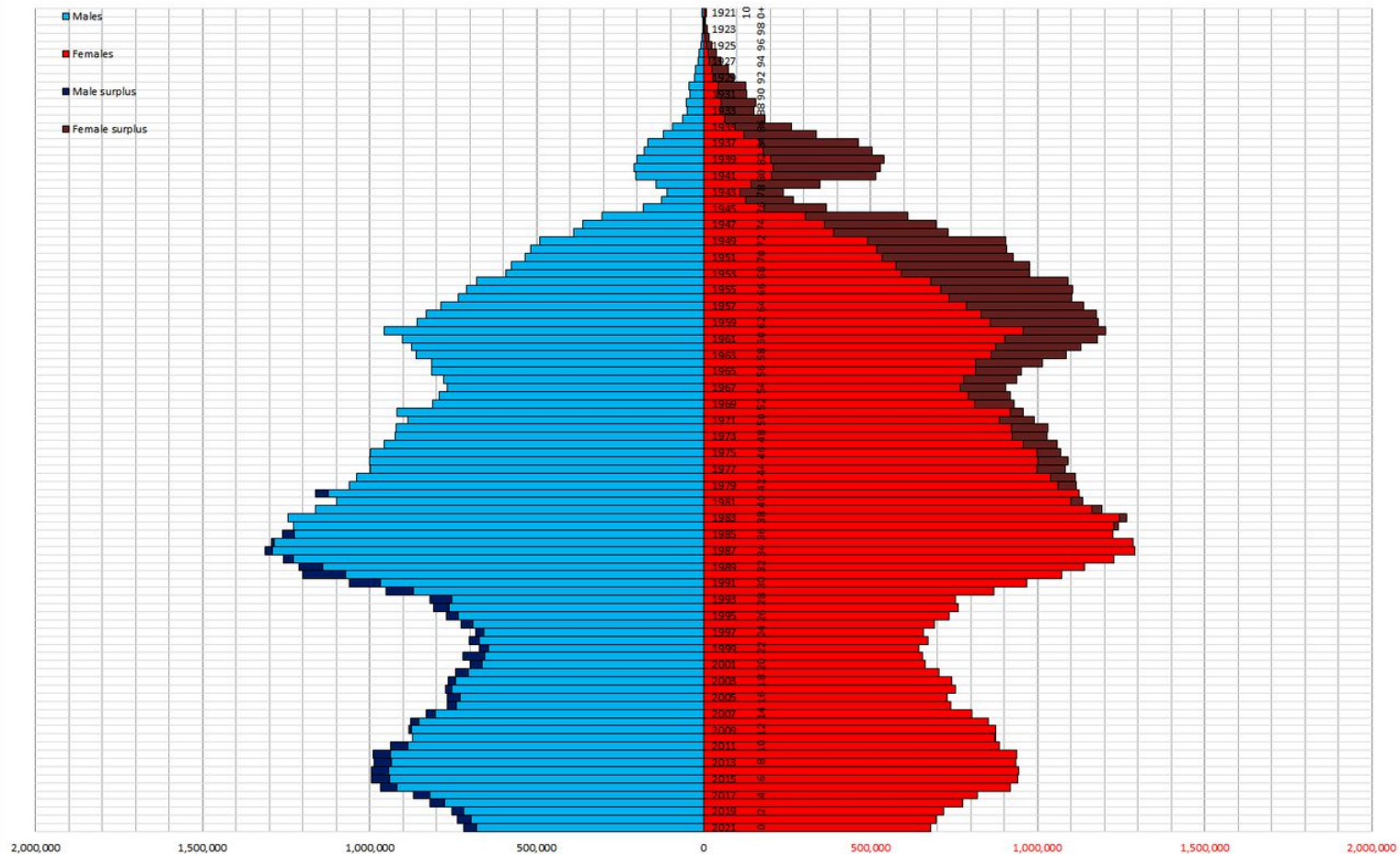
Ирис разноцветный
(*Iris versicolor*)



Ирис виргинский
(*Iris virginica*)

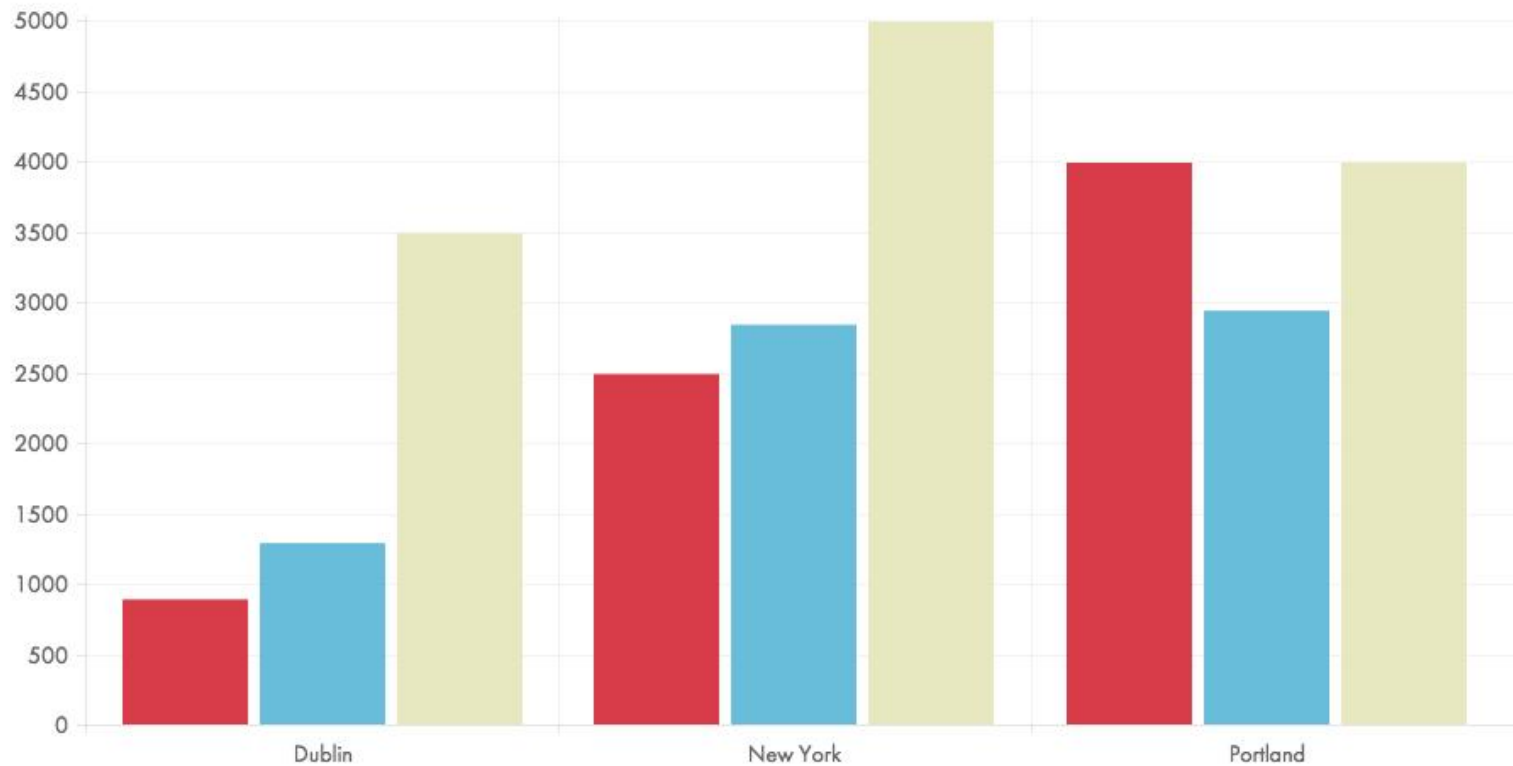


Russia Age Sex Structure 01.01.2022

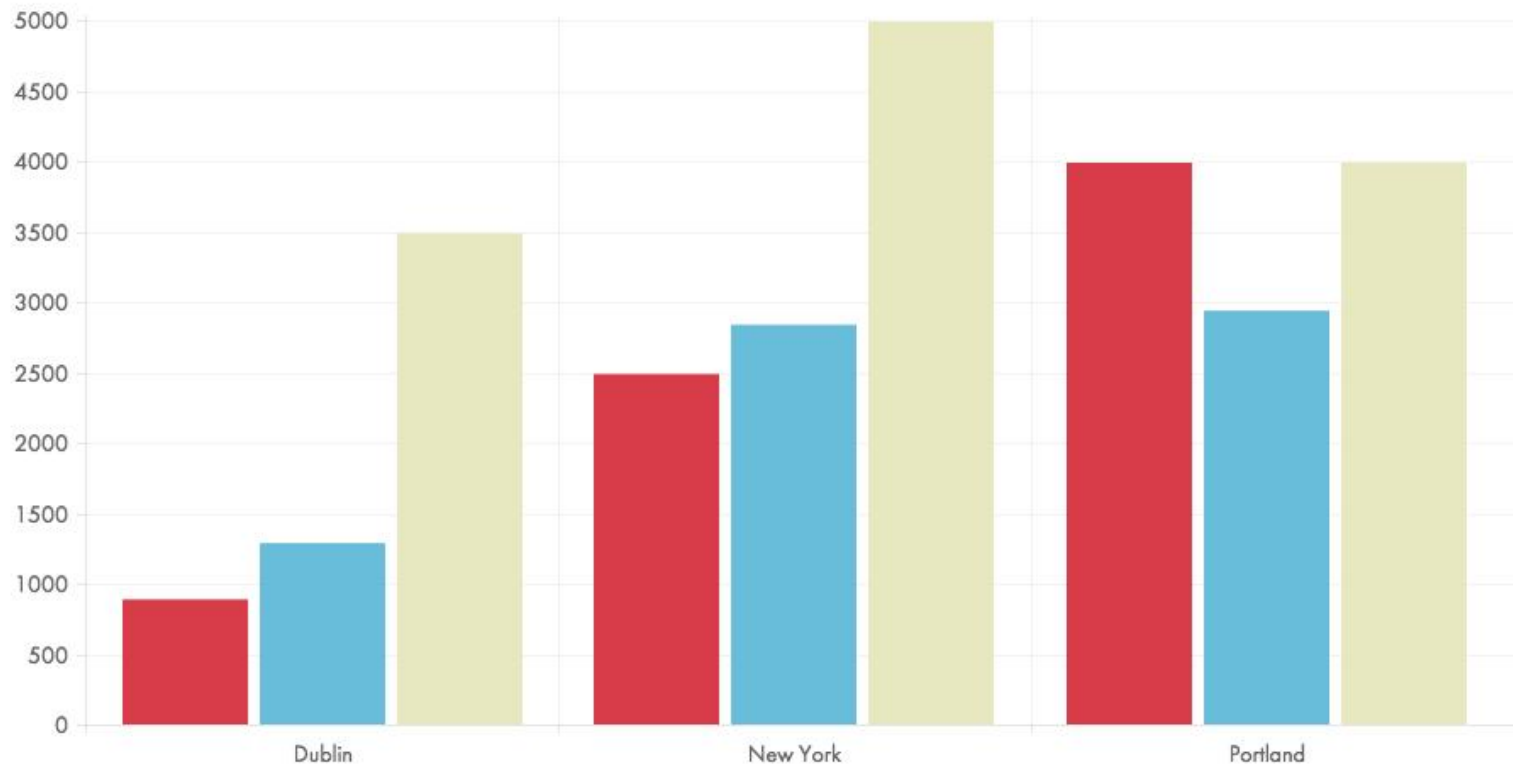


Bar chart

TICKET SALES BY LOCATION (JANUARY TO MARCH)

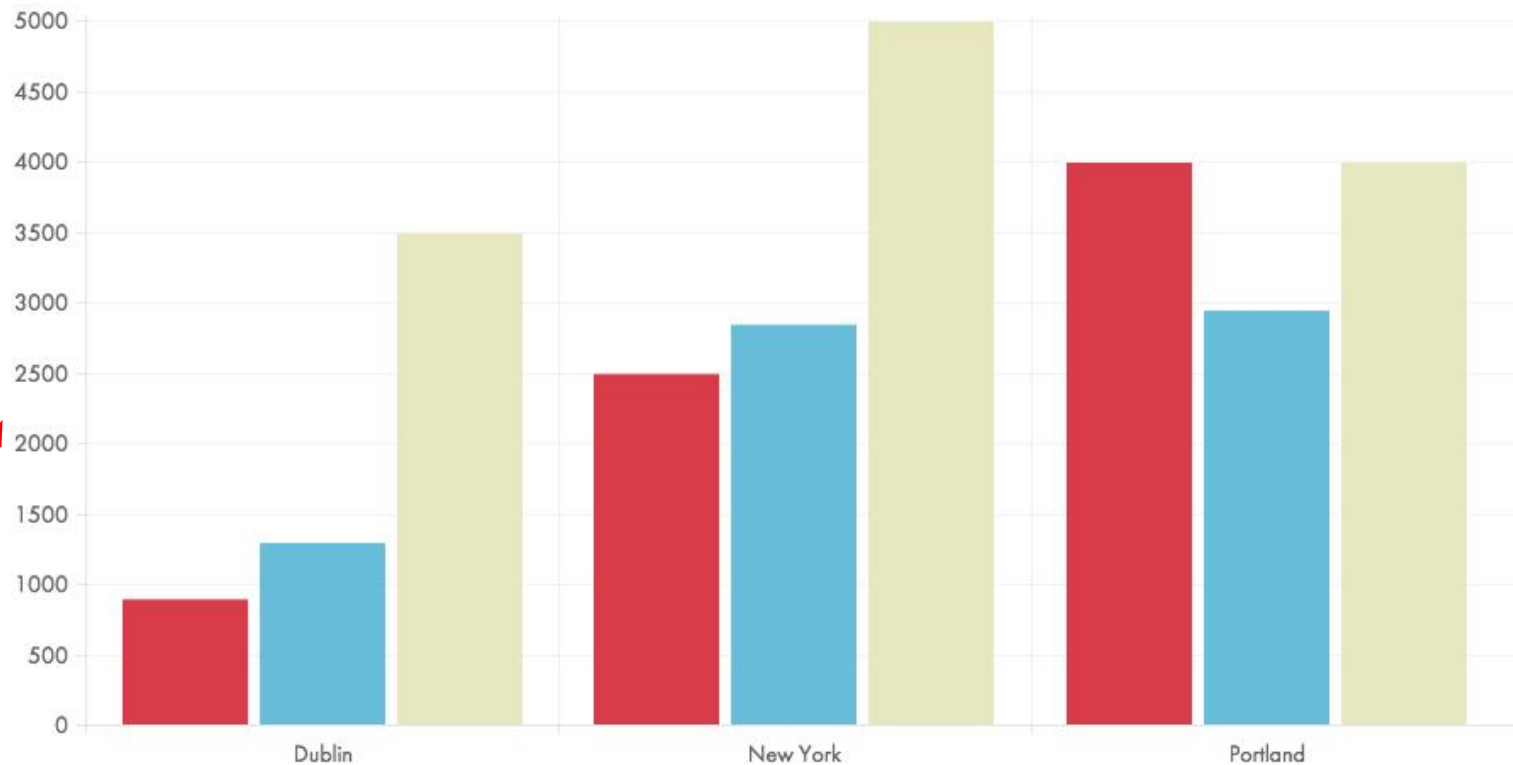


TICKET SALES BY LOCATION (JANUARY TO MARCH)



← Category

TICKET SALES BY LOCATION (JANUARY TO MARCH)

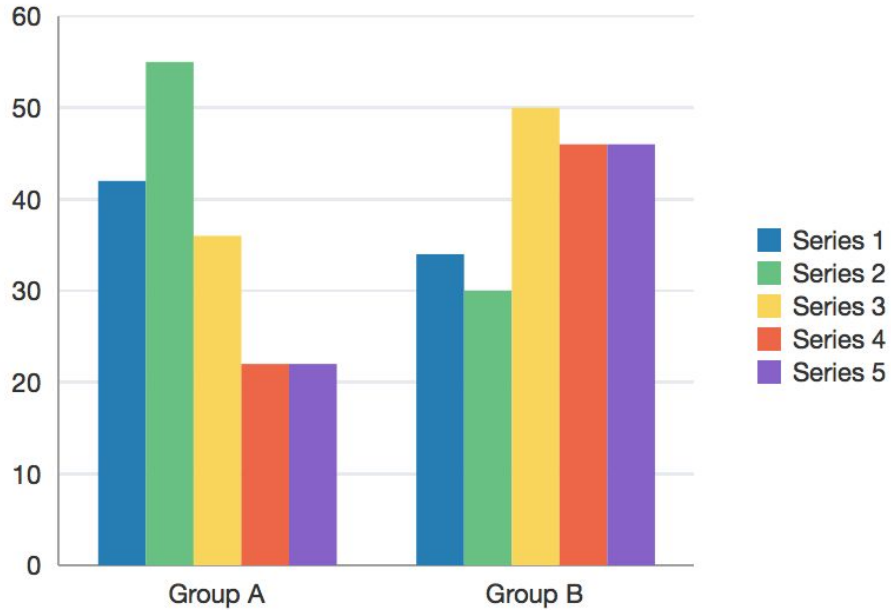


Value

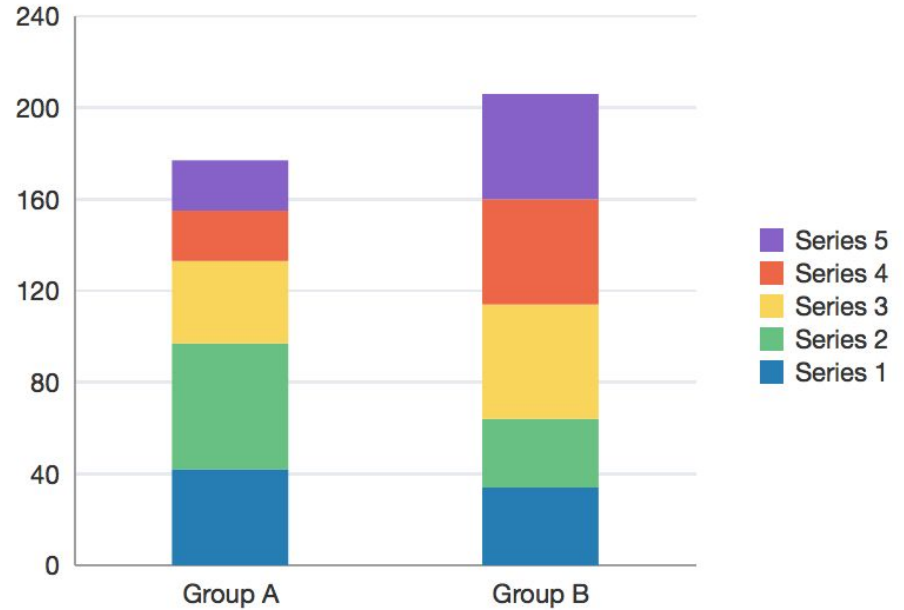
Category

Stacked bar chart

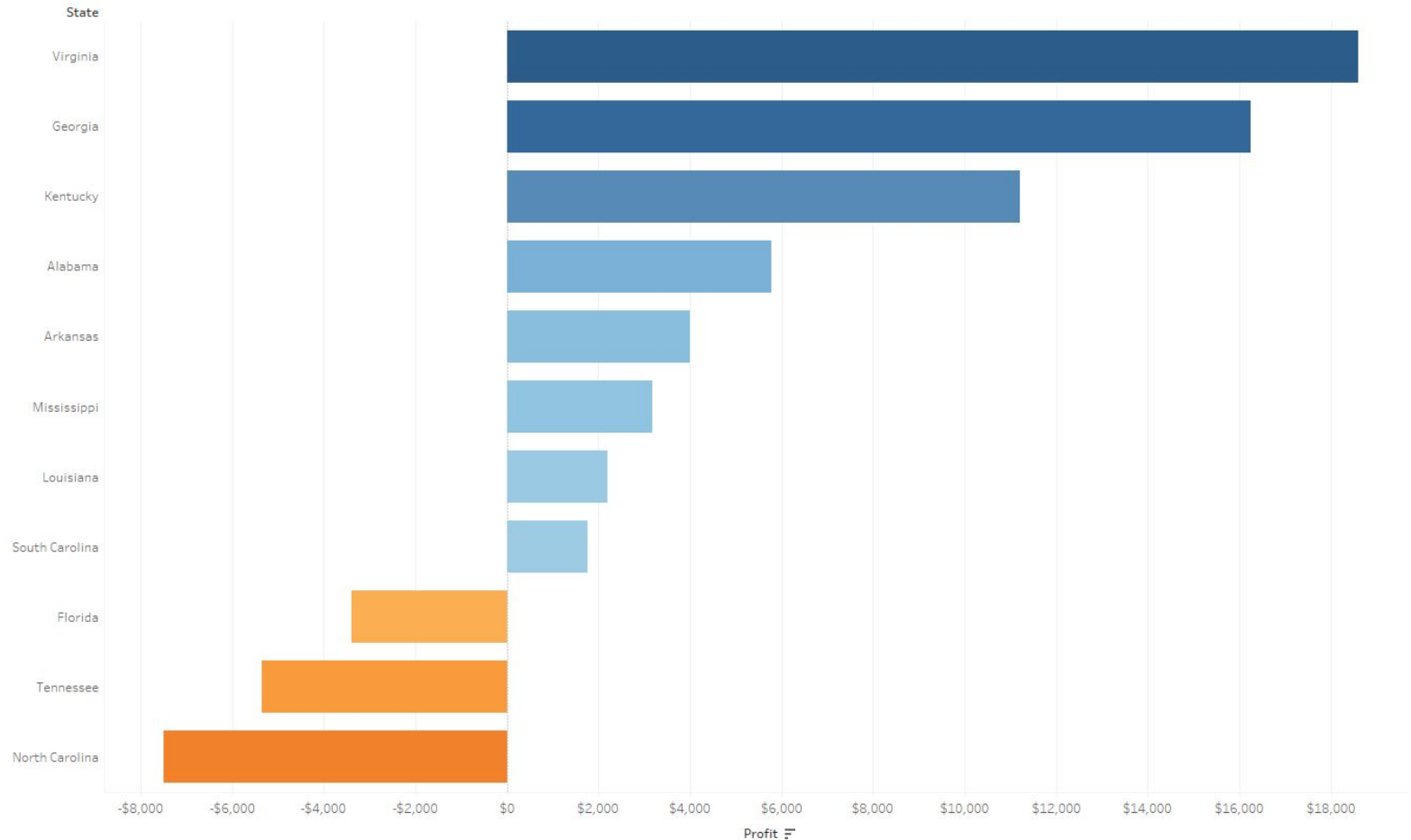
Bar Chart



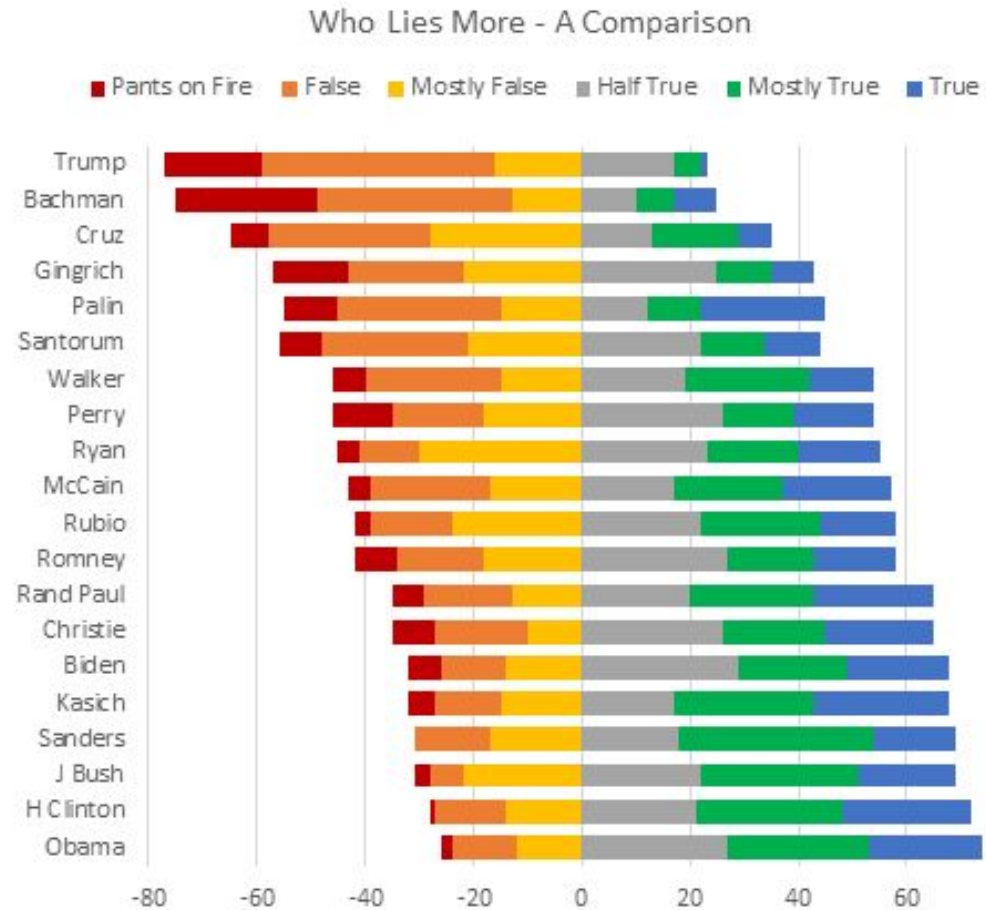
Stacked Bar Chart



Vertical bar chart



Vertical stacked bar chart



Box-and-whiskers diagram (box plot)



Ящик с усами

Материал из Википедии — свободной энциклопедии

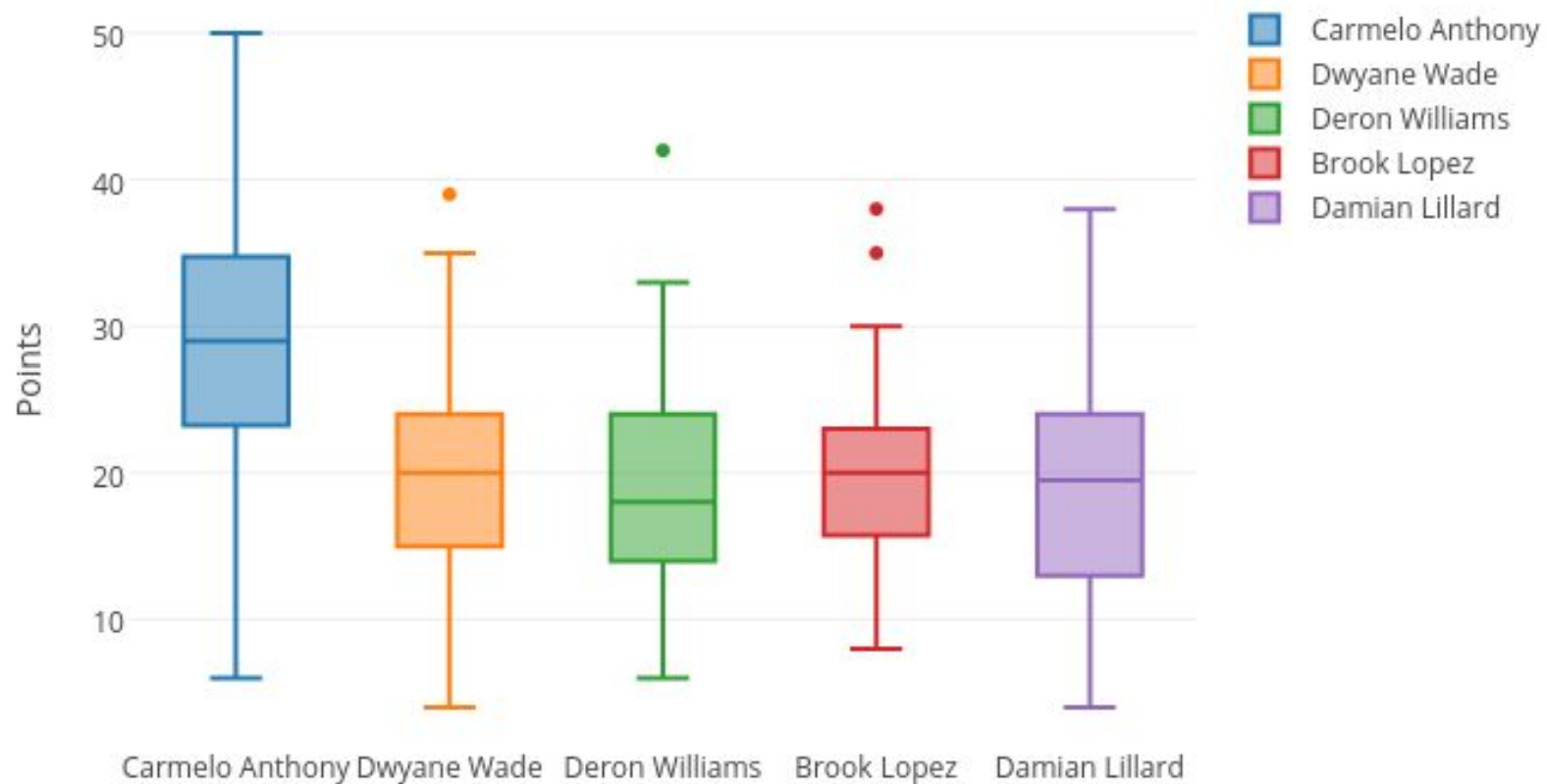
Не следует путать с [японскими свечами](#).

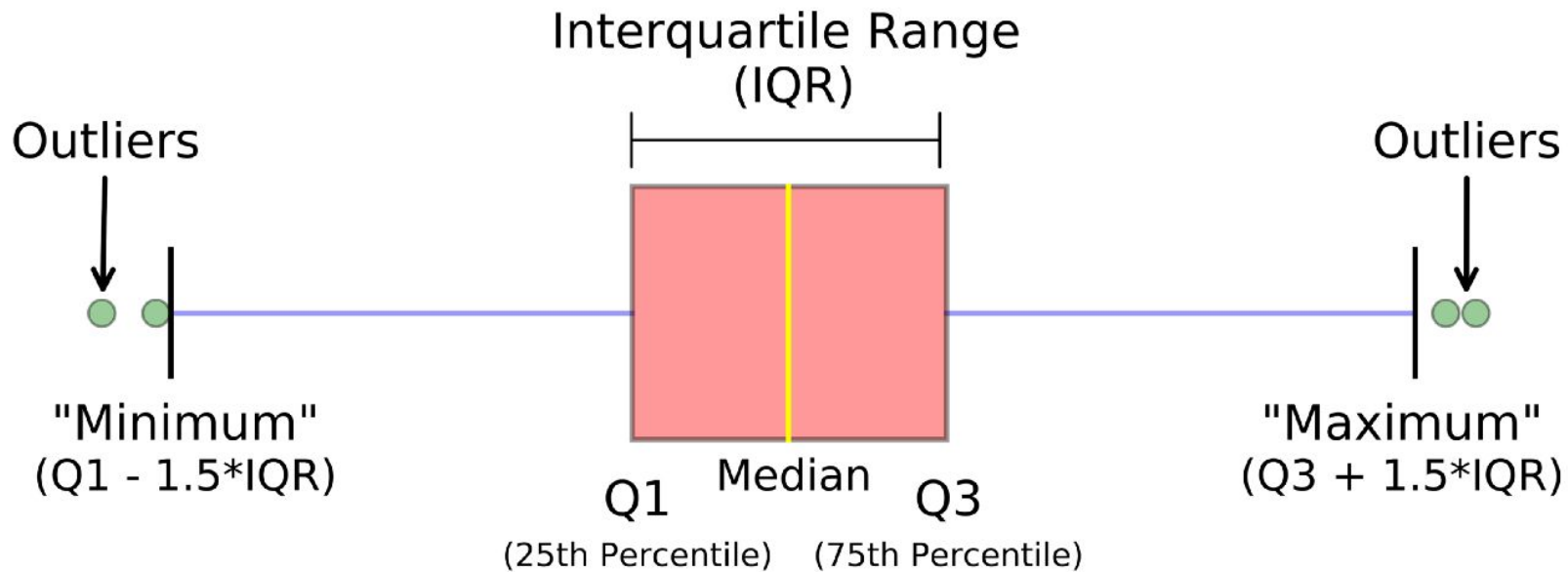
Ящик с усами, **диаграмма размаха** (англ. *box-and-whiskers diagram or plot, box plot*, вероятностей).

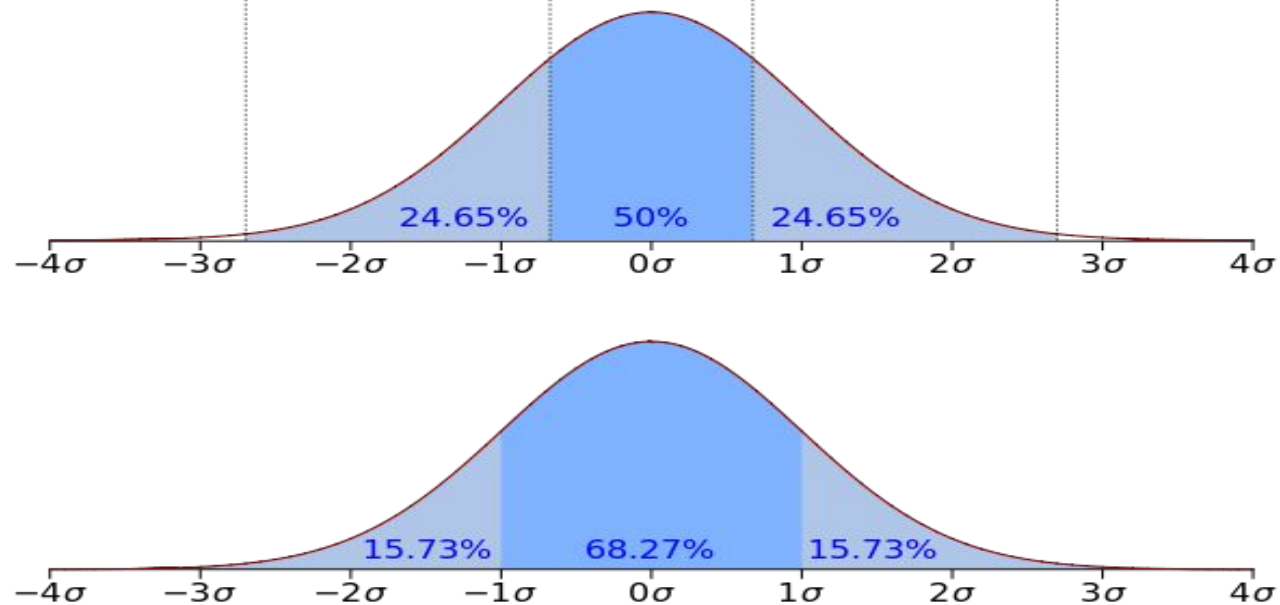
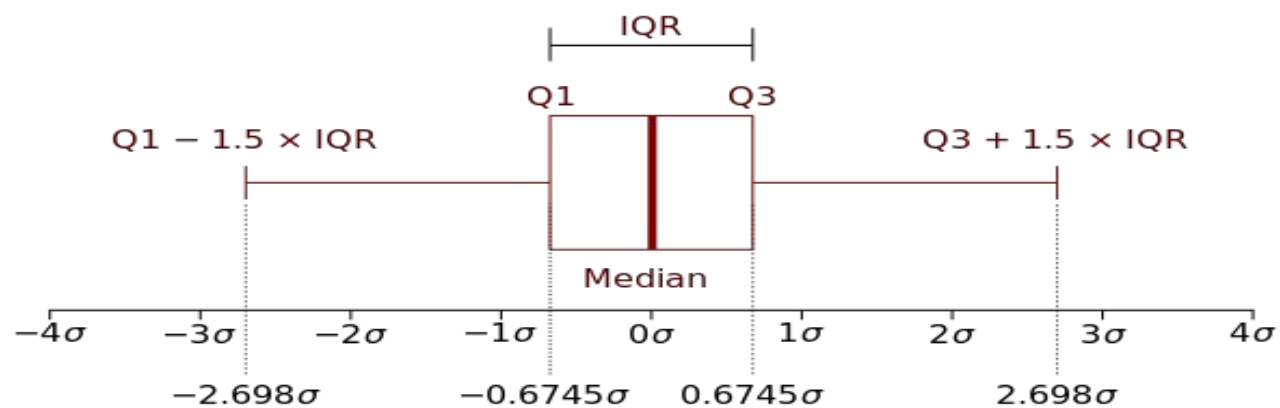
Такой вид диаграммы в удобной форме показывает медиану (или, если нужно, сред

[Wikipedia](#)

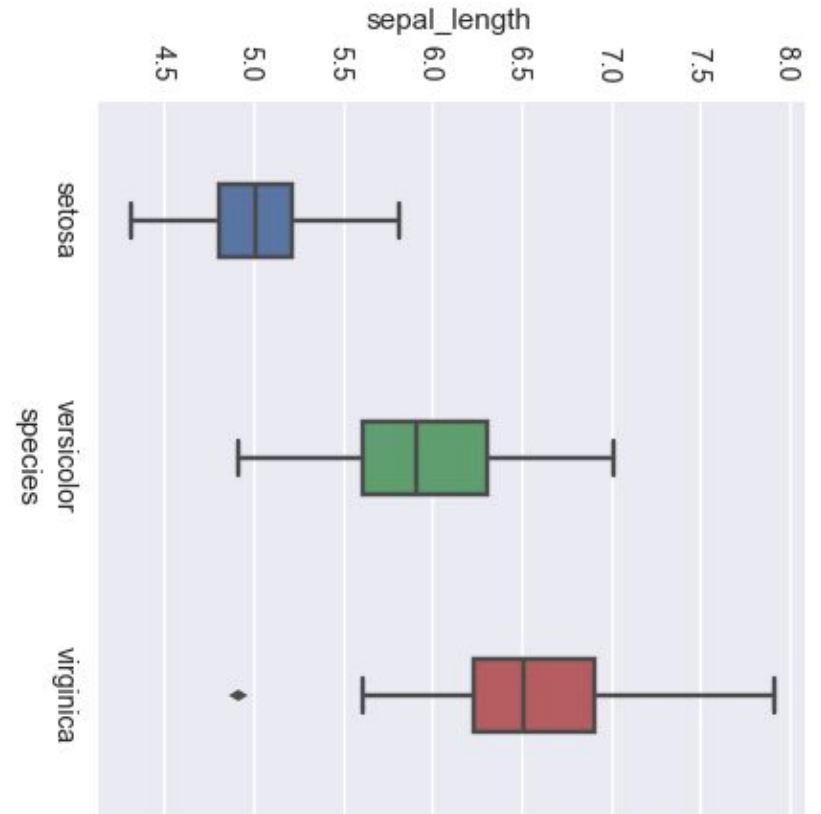
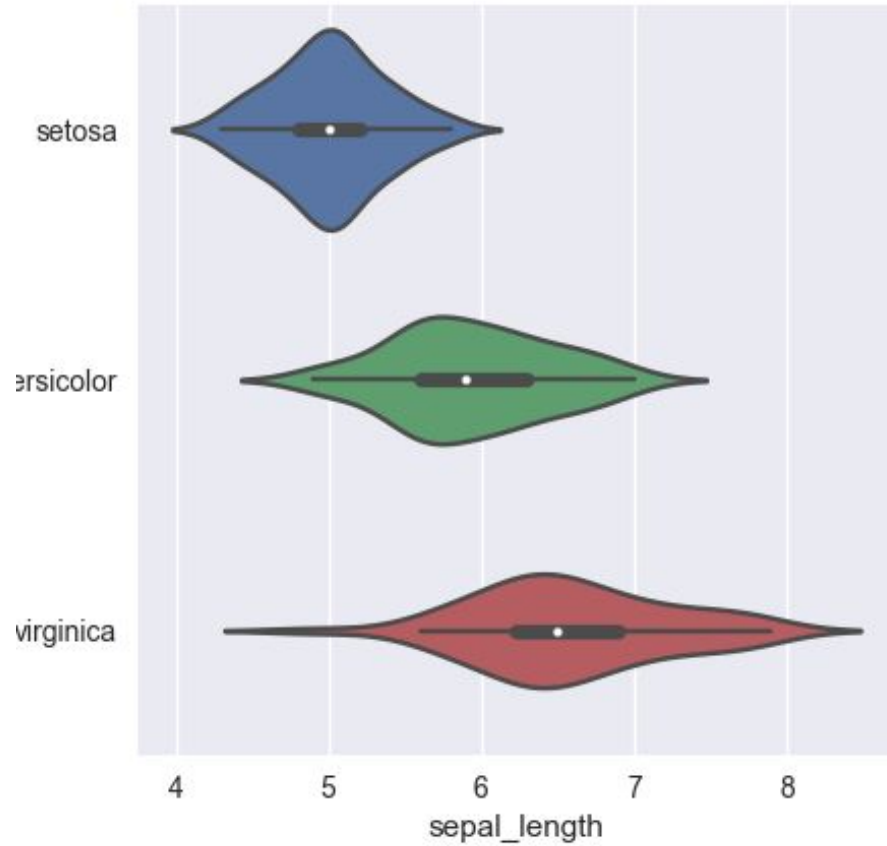
Points Scored Per NBA Game

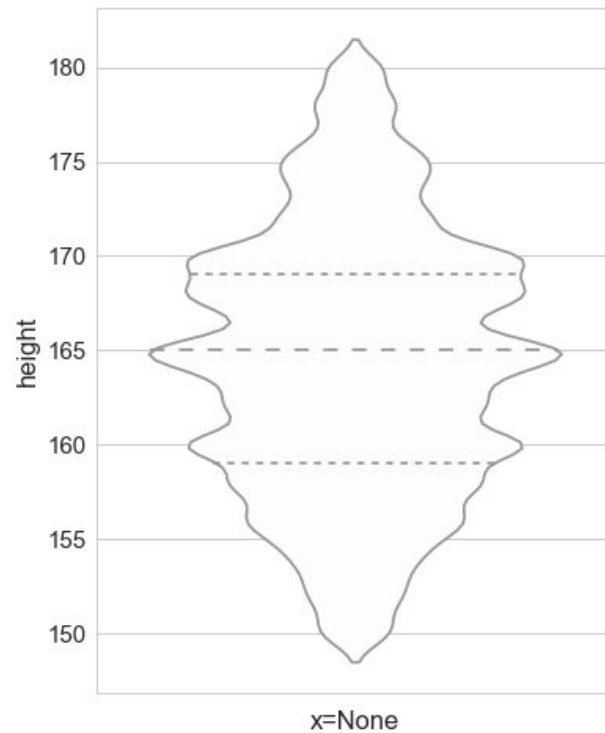
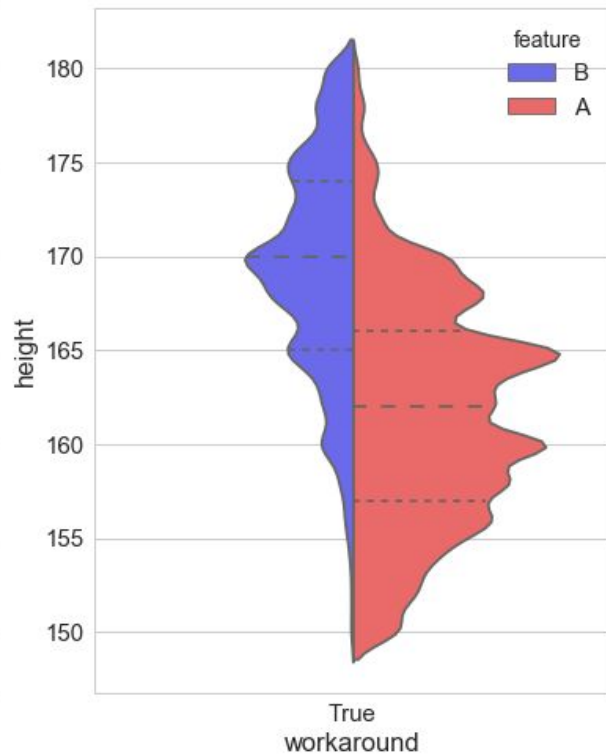
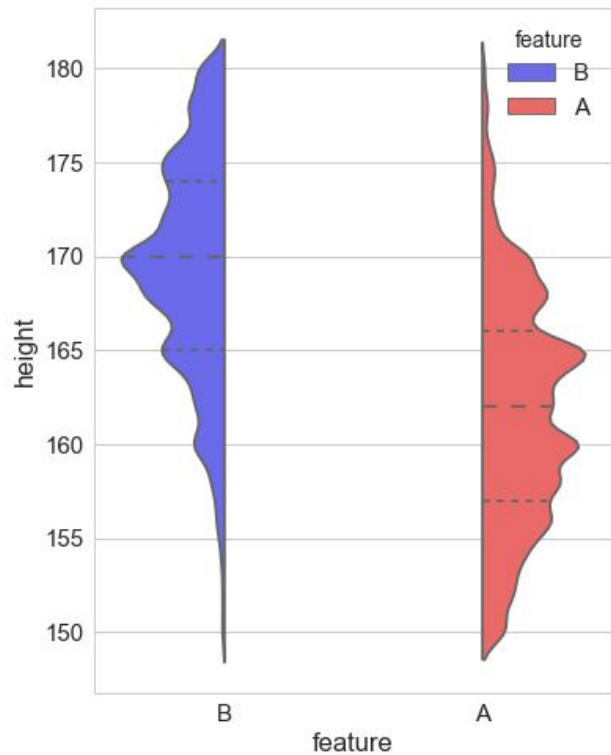






Violin plot





Pie chart

Пицца – это круговая диаграмма,
показывающая, сколько
у тебя осталось пиццы.

101

6



Пирог тоже

today at 8:42

2



Умно

today at 8:46

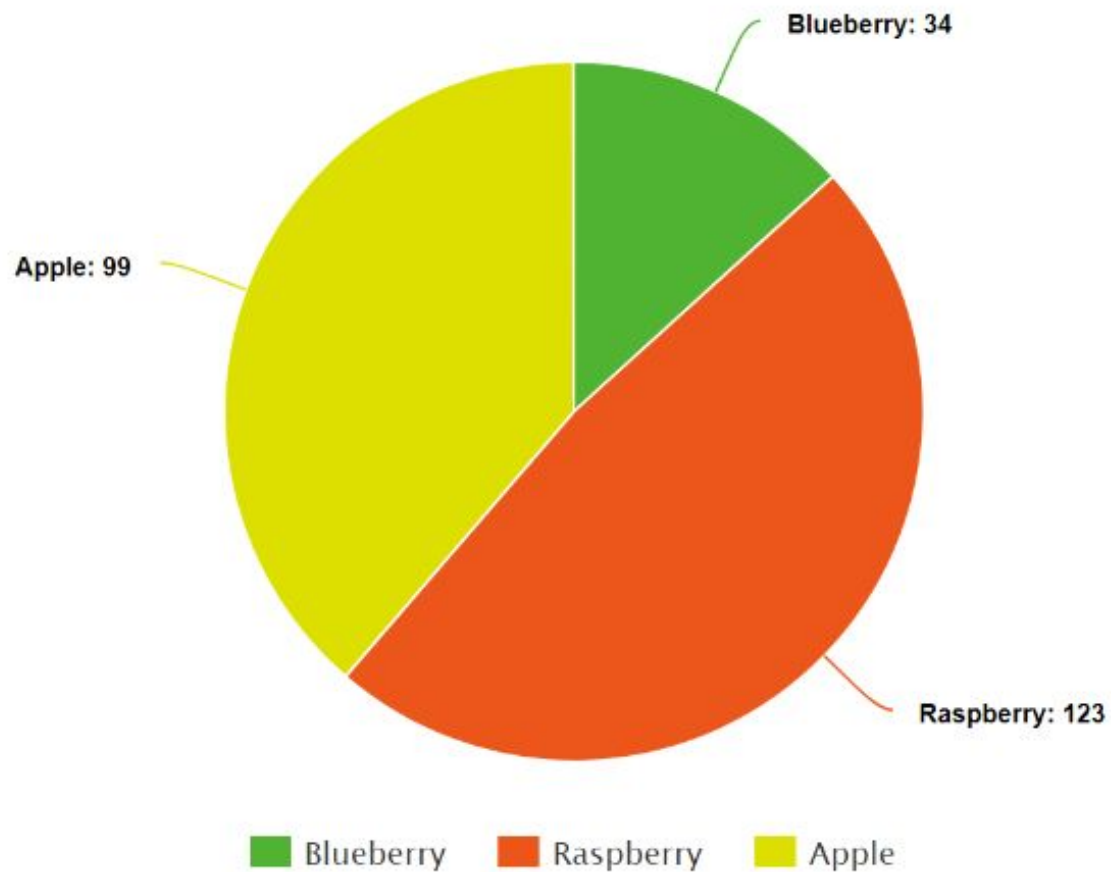


Вадим, пирог не показывает, сколько у тебя осталось пиццы. Так что - нет,
не то-же.

today at 8:46 to Vadim

24

Pie of Pies

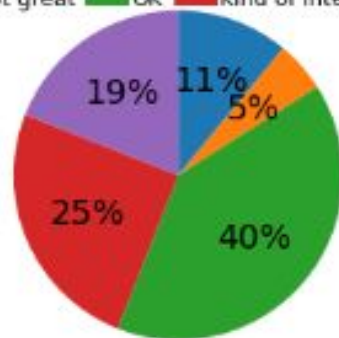


Pie charts are the
worst

Survey Results

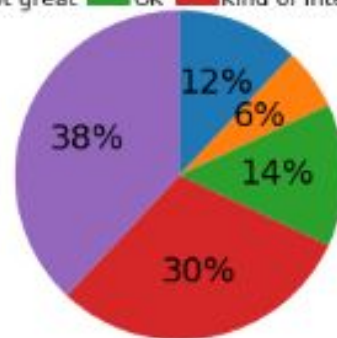
PRE: How do you feel about doing science?

■ Bored ■ Not great ■ OK ■ Kind of interested ■ Excited



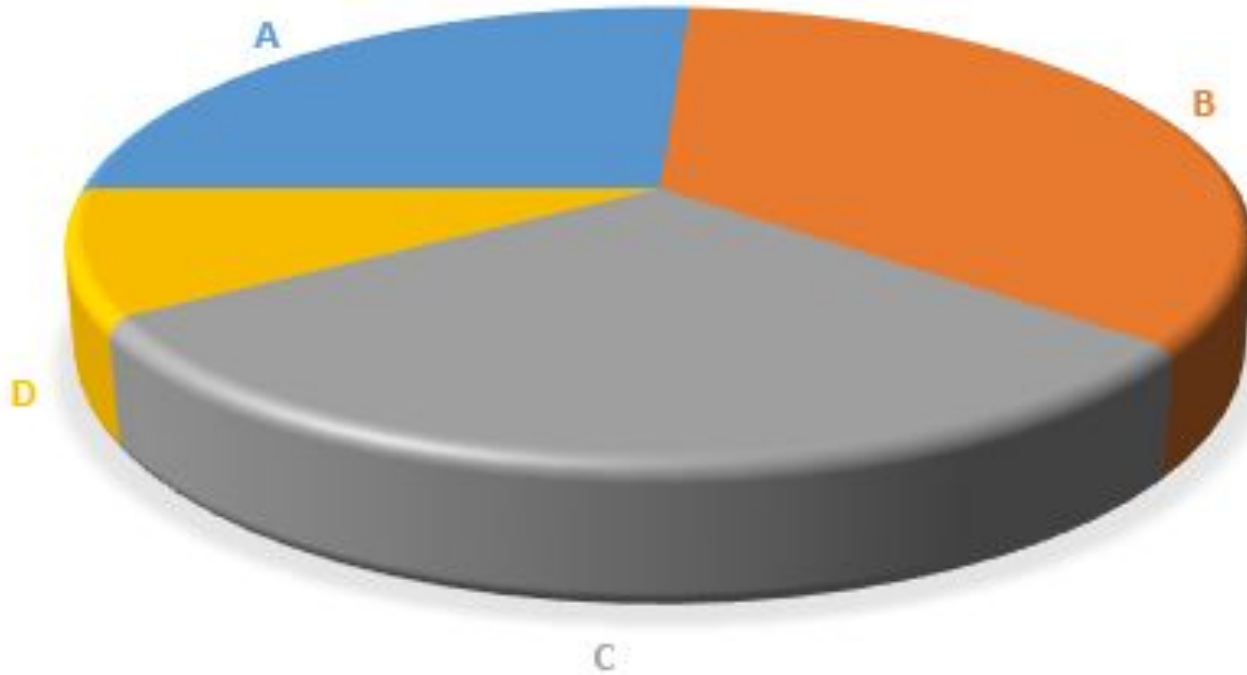
POST: How do you feel about doing science?

■ Bored ■ Not great ■ OK ■ Kind of interested ■ Excited

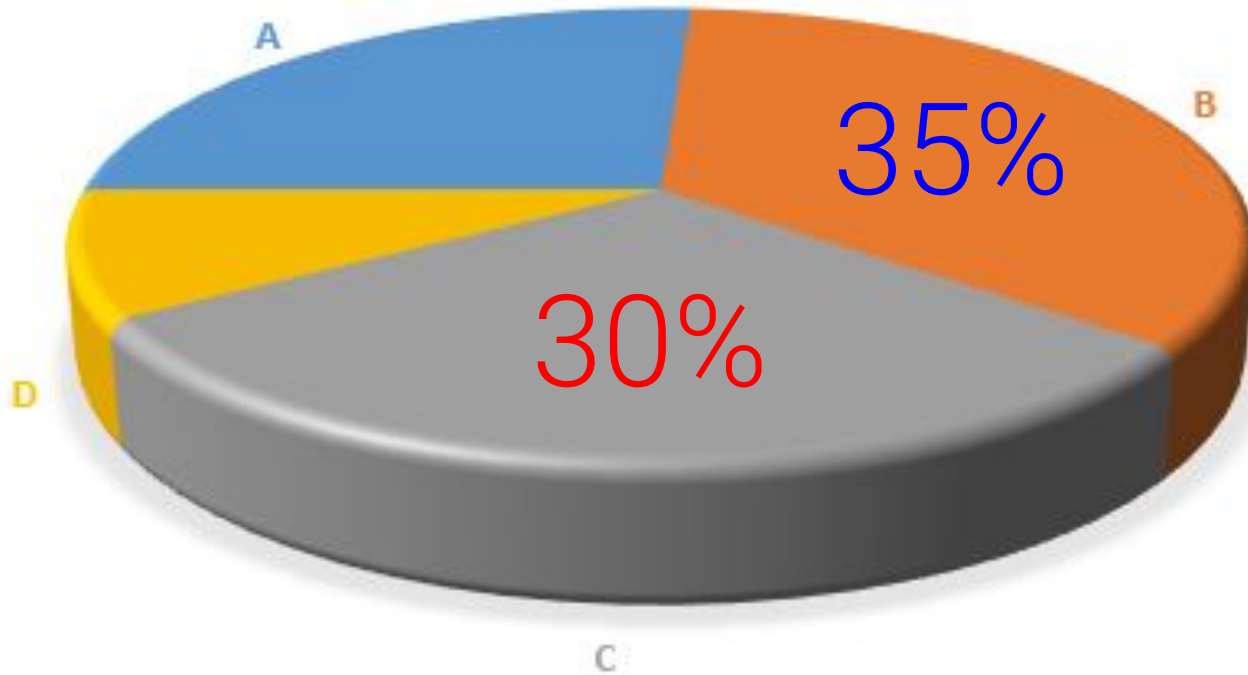


3D Pie charts
are evil

Which is bigger C or B?

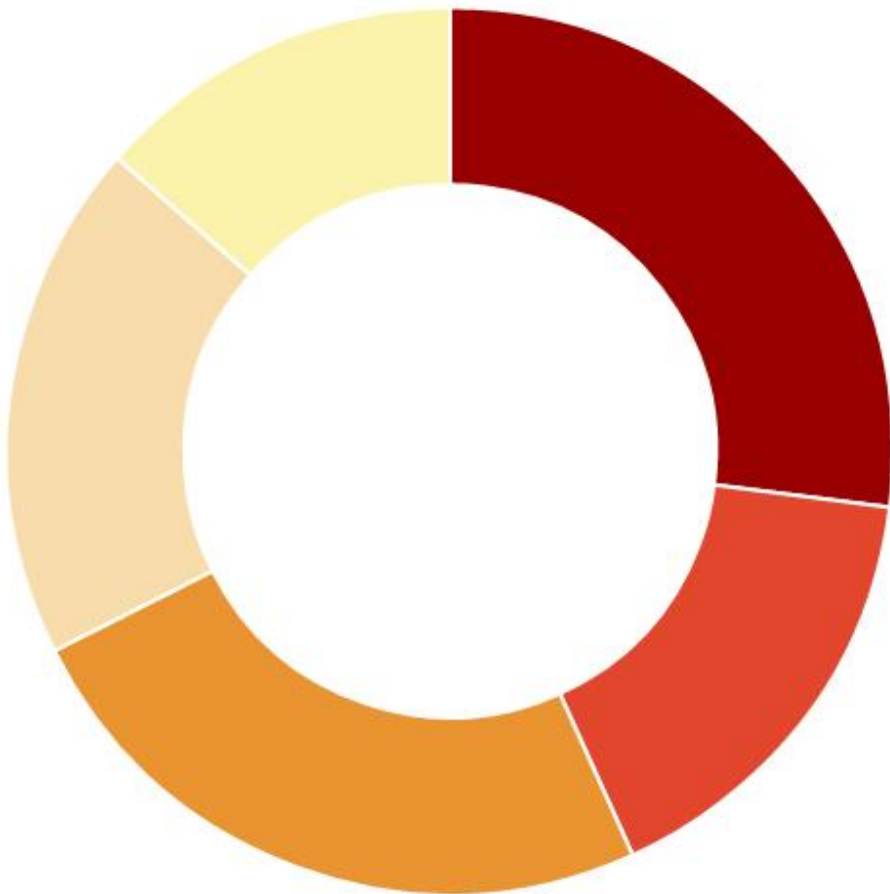


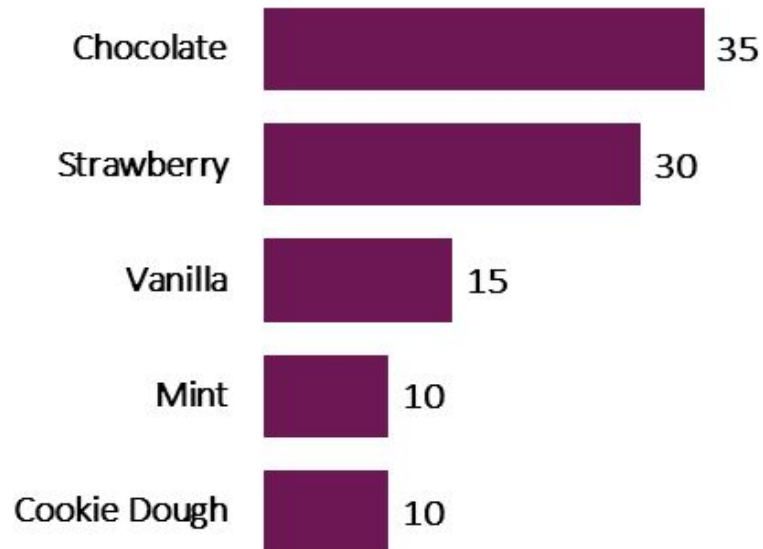
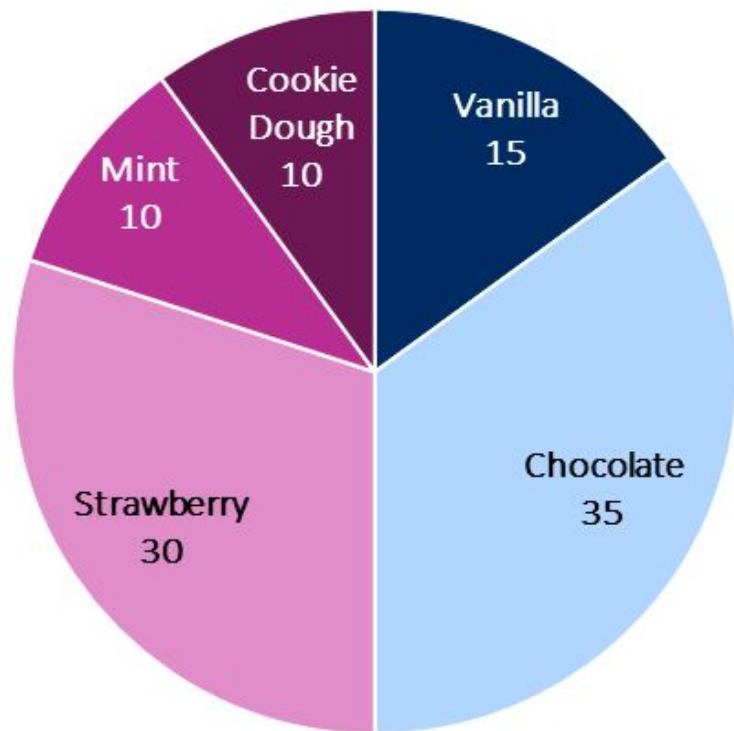
Which is bigger C or B?



FAVORITE PIES

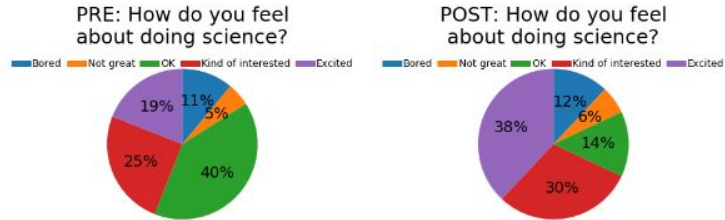
- Apple Pie
- Pumpkin Pie
- Pecan Pie
- Cherry Pie
- Other





Было:

Survey Results

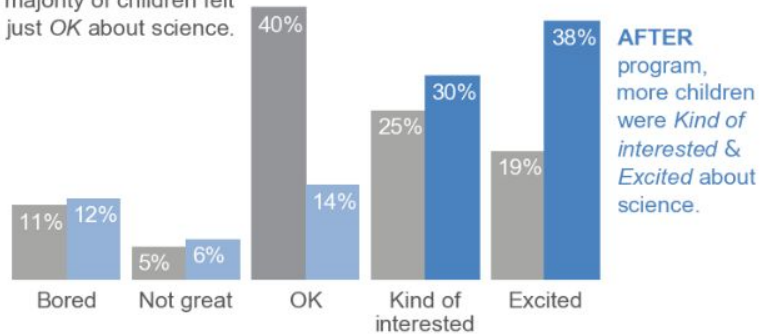


Стало:

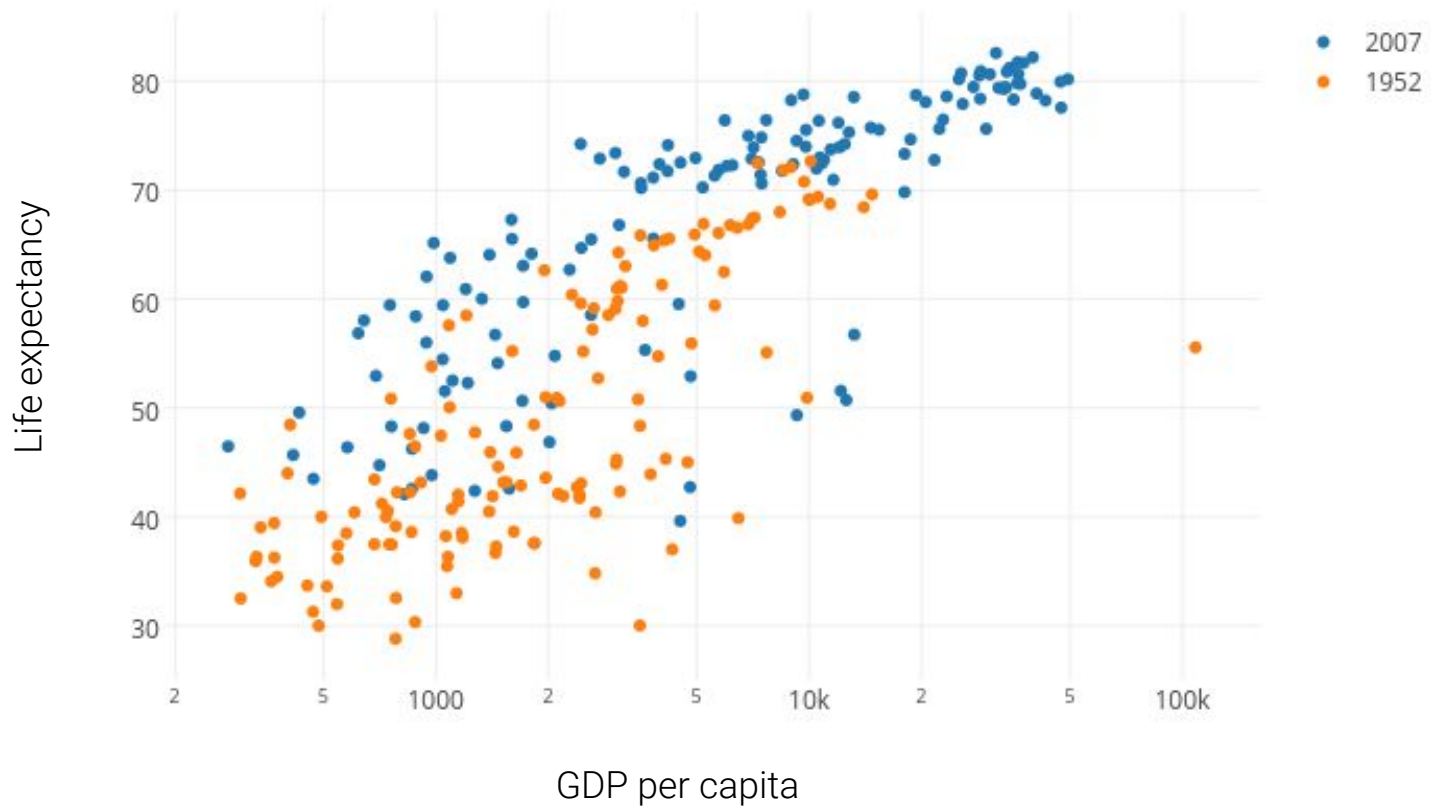
Pilot program was a success

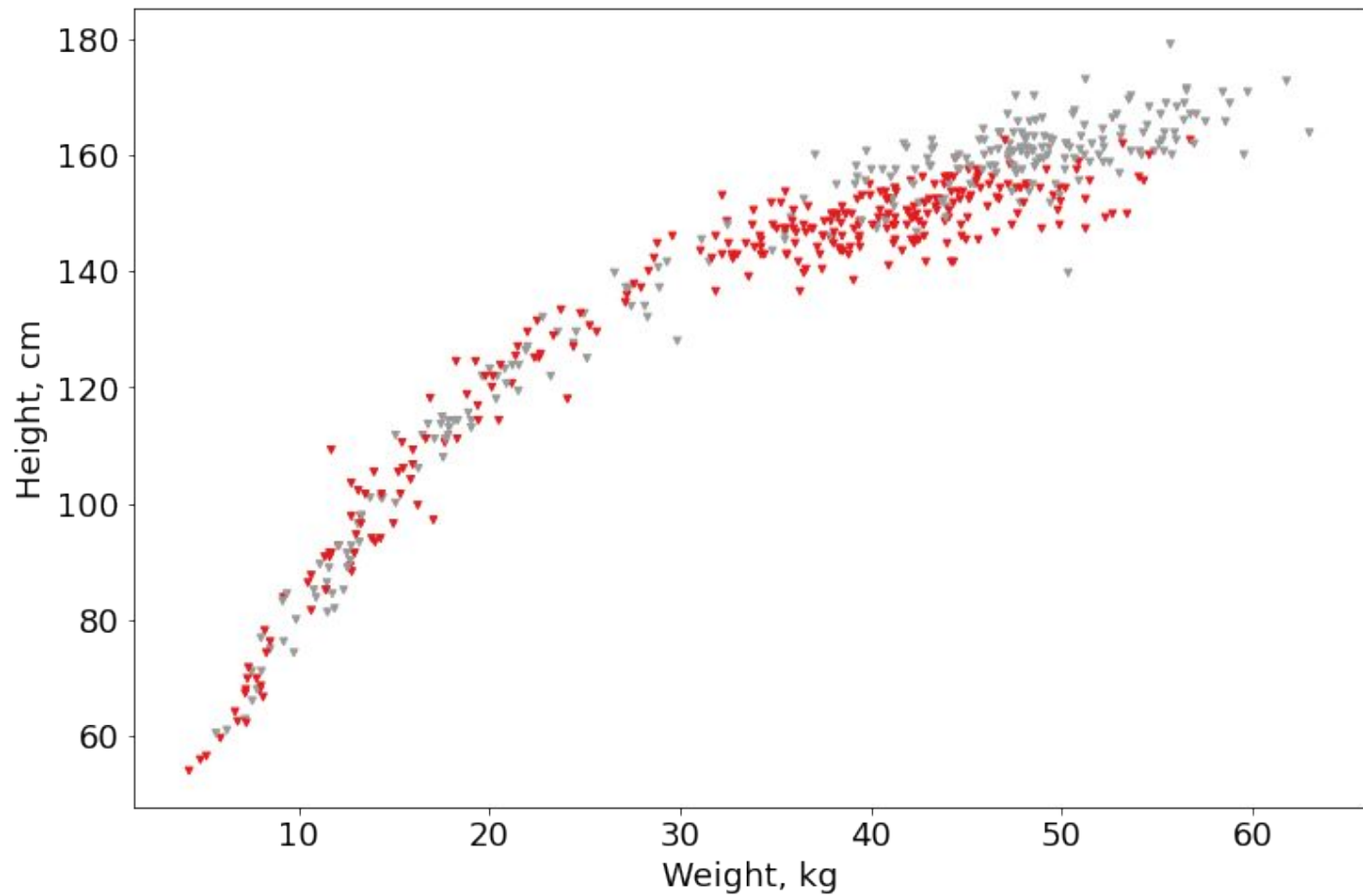
How do you feel about science?

BEFORE program, the majority of children felt just OK about science.



Scatter plot



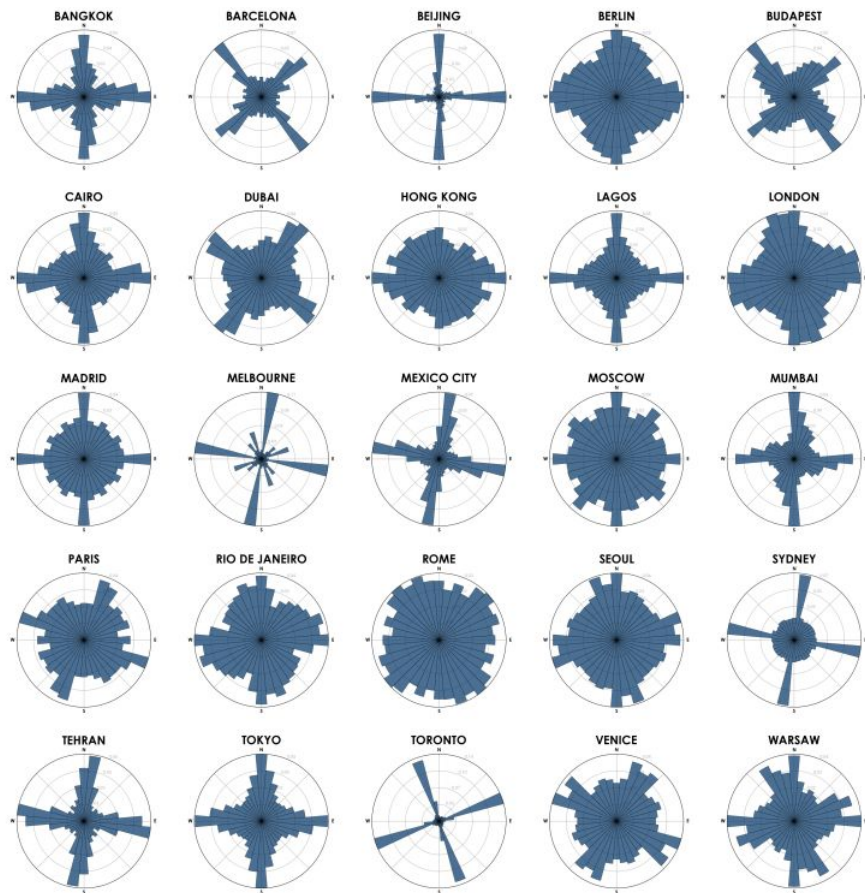


What else?

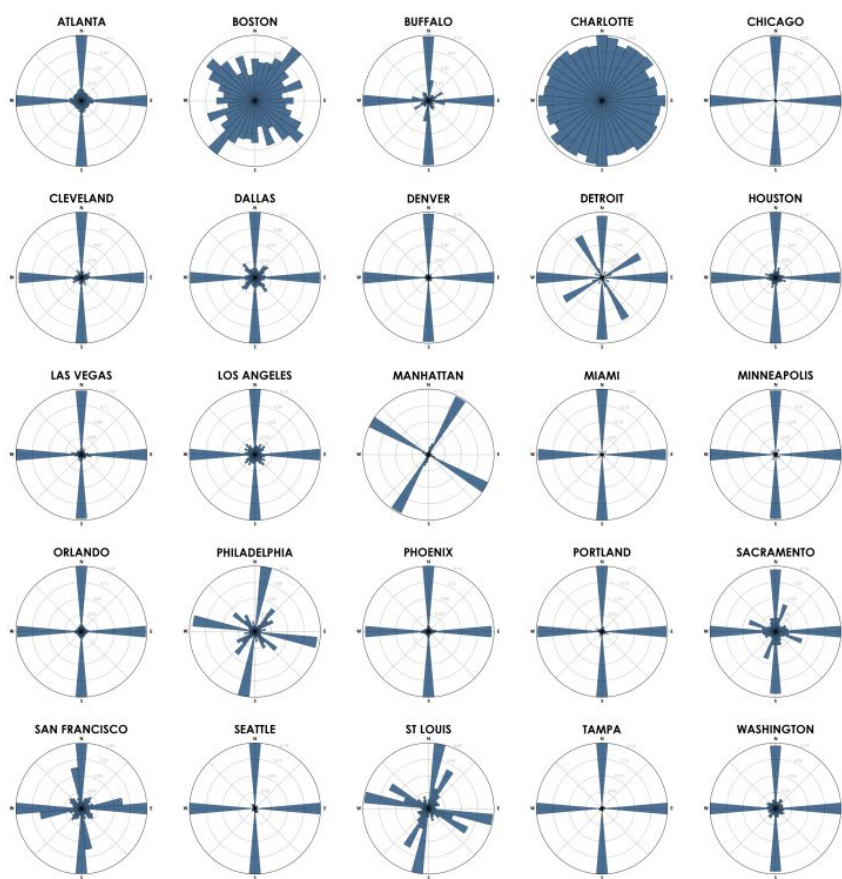
Other types of visualizations

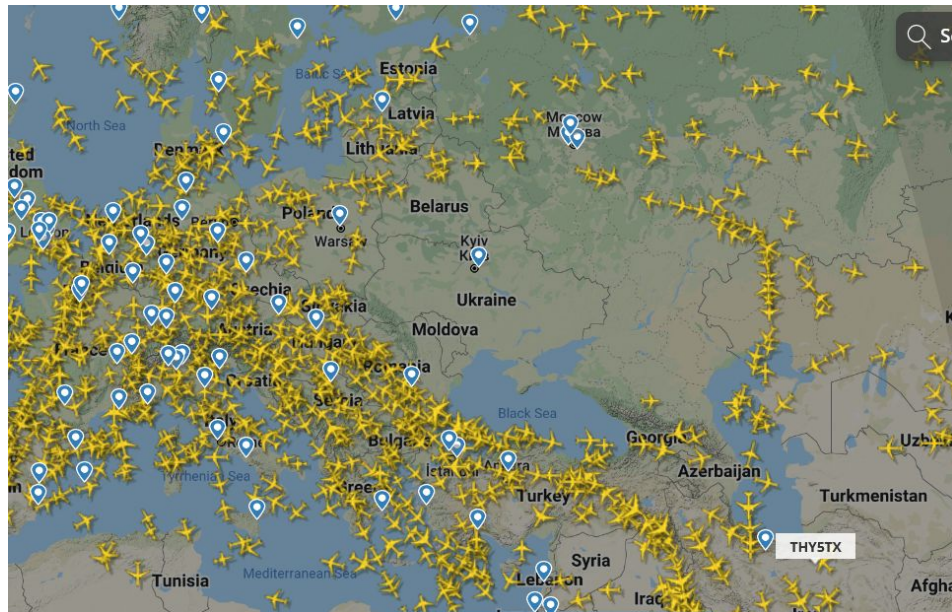
1. Slope chart
2. Venn diagrams
3. Geo maps
4. Bubble chart
5. Tag cloud
6. Stacked graphs
7. Stem and Leaf plot
8. Networks (social, information, etc)
9. ...

City Street Network Orientation

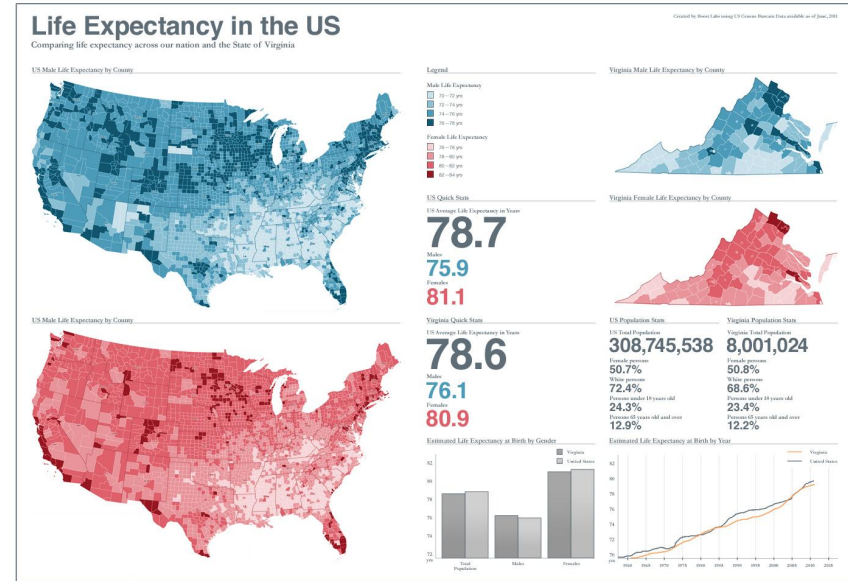


City Street Network Orientation



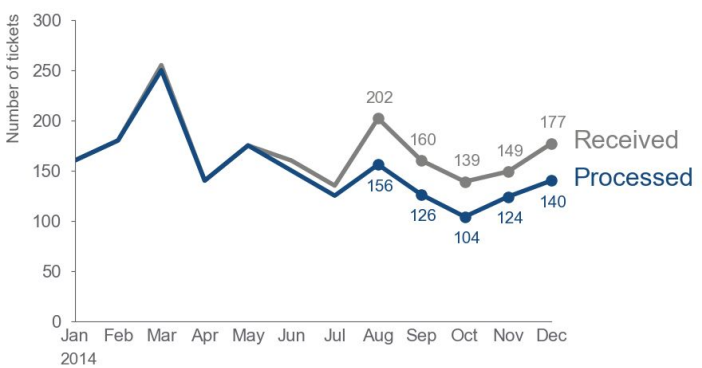


<https://www.flightradar24.com>



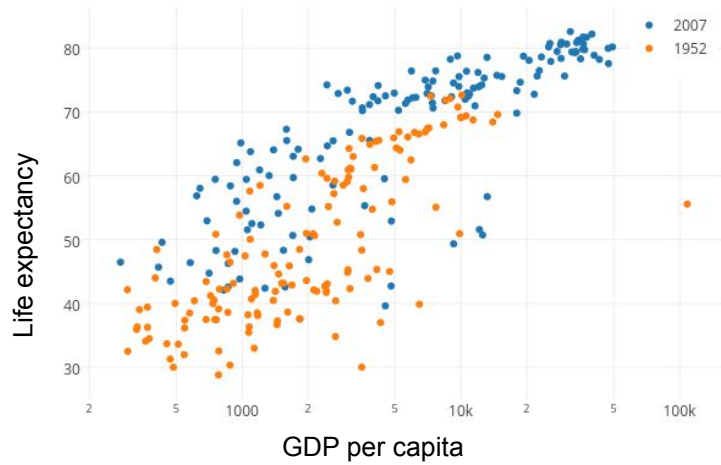
<https://blogs.stockton.edu/datavis/2017/11/27/geographical-visualizations/>

Ticket volume over time

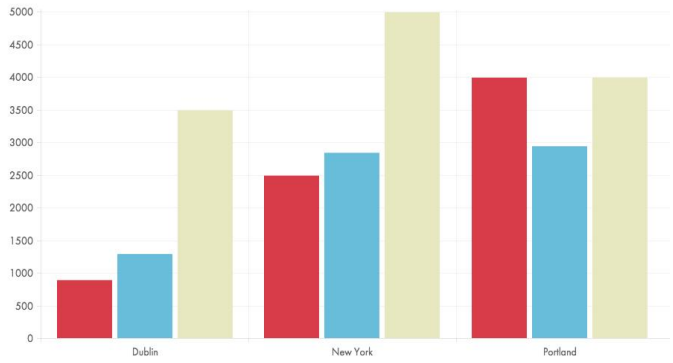


Data source: XYZ Dashboard, as of 12/31/2014

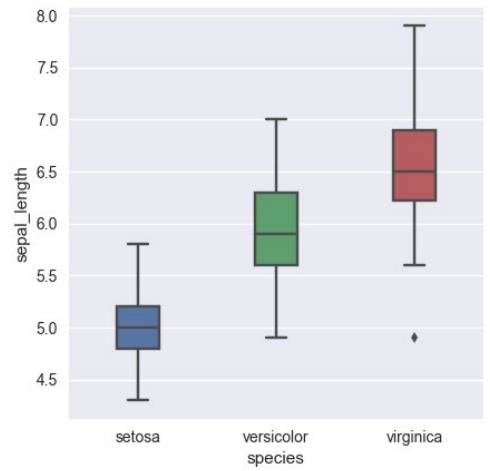
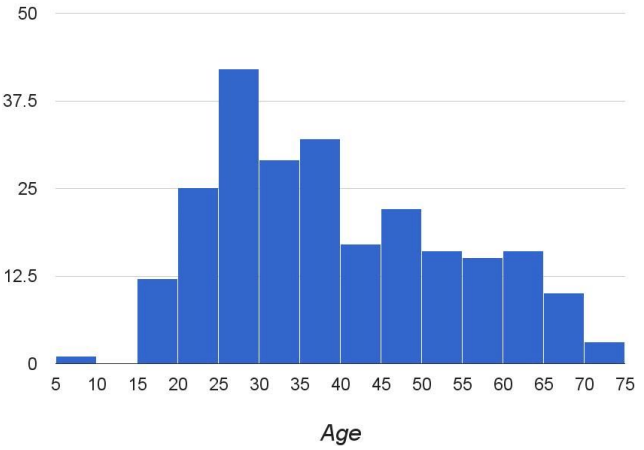
- Index chart
- Scatter plot
- Bar chart
- Histogram
- Box plot



TICKET SALES BY LOCATION (JANUARY TO MARCH)



Respondents' age distribution (n=240)



Orange <https://orangedatamining.com/>

KNIME <https://www.knime.com/>

Rapid Miner <https://rapidminer.com/>

Microsoft Azure <https://azure.microsoft.com/en-us/>