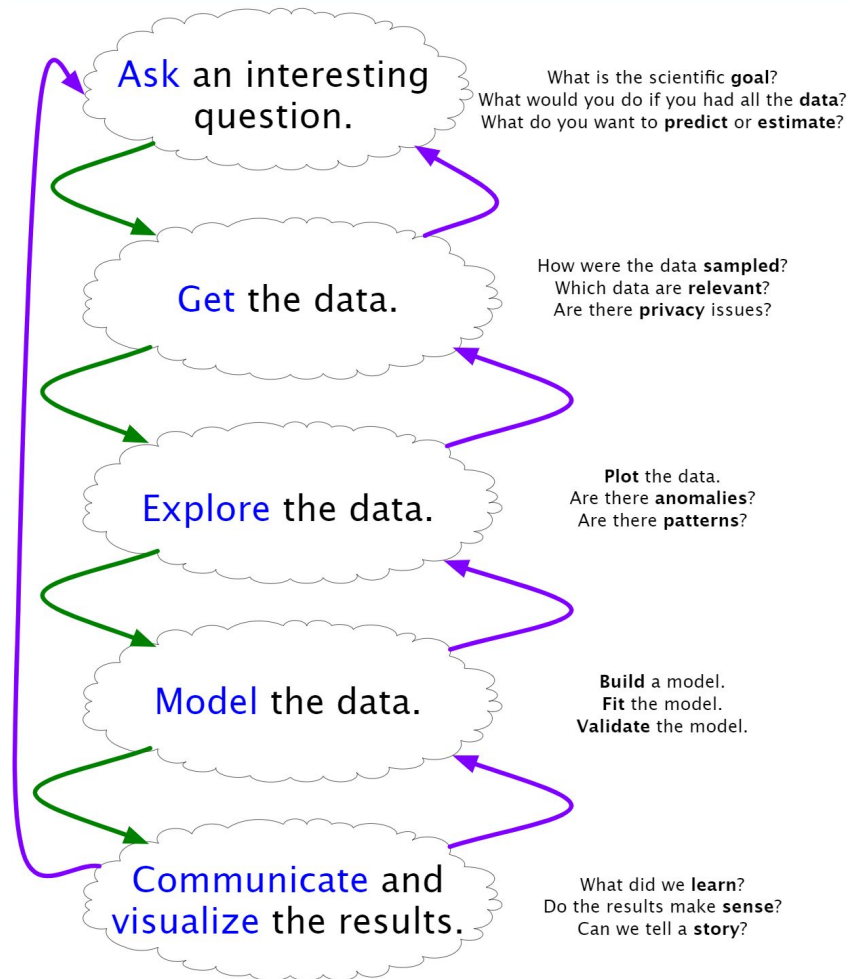


Practical Data Science

Week 2 - Lecture 2

**Exploratory Data Analysis & Effective
Visualizations**



Data exploration

Not always sure what we are looking for
(until we find it)



Example: Antibiotics

Will Burtin

Data

Genus, Species

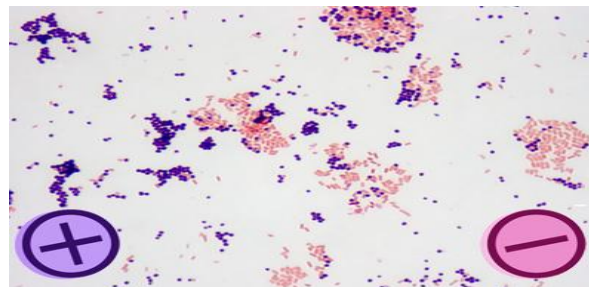
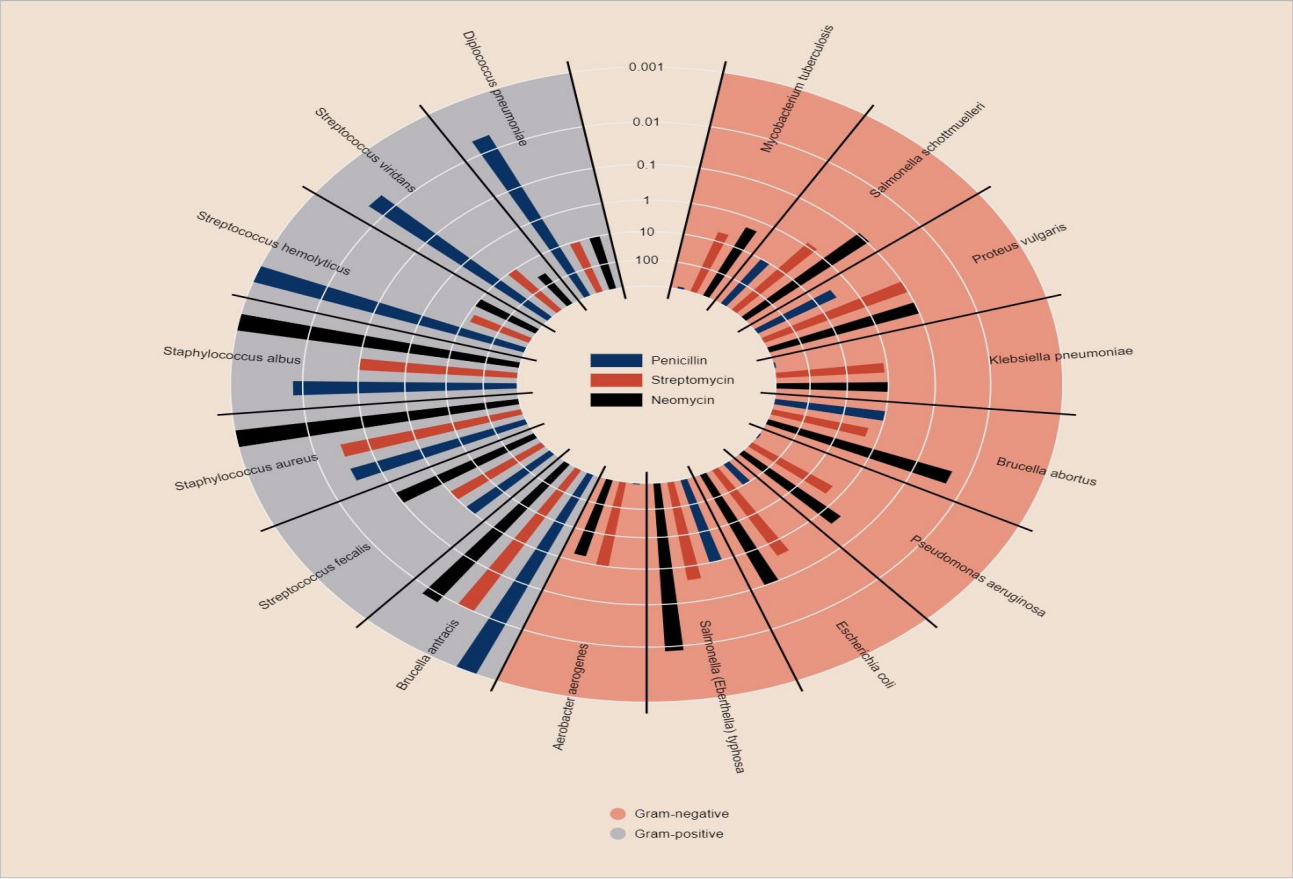


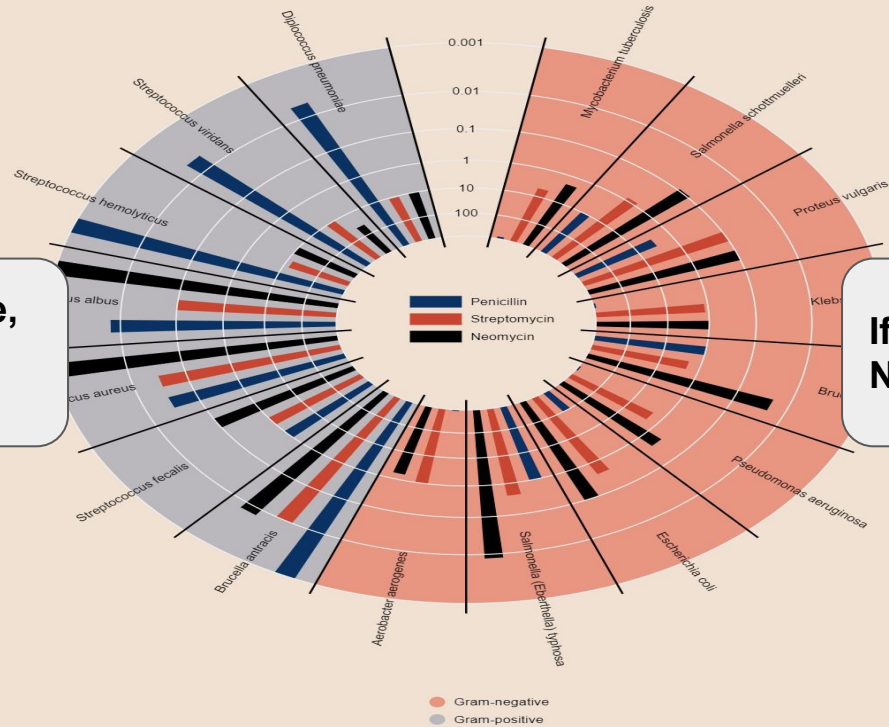
Table 1: Burtin's data.

Bacteria	Antibiotic			Gram Staining
	Penicillin	Streptomycin	Neomycin	
<i>Aerobacter aerogenes</i>	870	1	1.6	negative
<i>Brucella abortus</i>	1	2	0.02	negative
<i>Brucella anthracis</i>	0.001	0.01	0.007	positive
<i>Diplococcus pneumoniae</i>	0.005	11	10	positive
<i>Escherichia coli</i>	100	0.4	0.1	negative
<i>Klebsiella pneumoniae</i>	850	1.2	1	negative
<i>Mycobacterium tuberculosis</i>	800	5	2	negative
<i>Proteus vulgaris</i>	3	0.1	0.1	negative
<i>Pseudomonas aeruginosa</i>	850	2	0.4	negative
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	negative
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	negative
<i>Staphylococcus albus</i>	0.007	0.1	0.001	positive
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	positive
<i>Streptococcus fecalis</i>	1	1	0.1	positive
<i>Streptococcus hemolyticus</i>	0.001	14	10	positive
<i>Streptococcus viridans</i>	0.005	10	40	positive

Burtin's Antibiotics



Burtin's Antibiotics



If bacteria is gram positive,
Penicillin & Neomycin are
most effective

If bacteria is gram negative,
Neomycin is most effective.

Exploratory Data Analysis

“The greatest value of a picture is when it forces us to notice what we never expected to see.”



John Tukey

Visualization

To convey information through
graphical representations of data

seaborn

0.9.0

Gallery

Tutorial

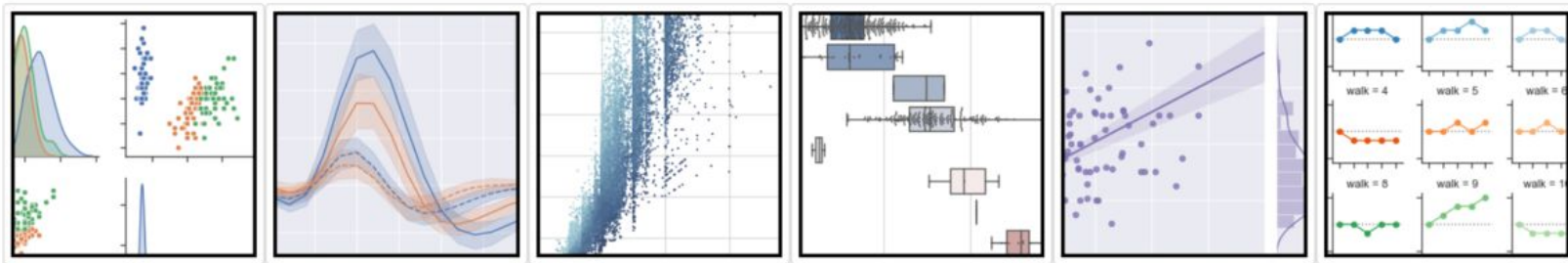
API

Site ▼

Page ▼

Search

seaborn: statistical data visualization



Visualization Goals

Communicate (Explanatory)

- Present data and ideas

- Explain and inform

- Provide evidence and support

- Influence and Persuade

Analyze (Exploratory)

- Explore the data

- Assess a situation

- Determine how to proceed

- Decide what to do

Effective visualizations

Not Effective...

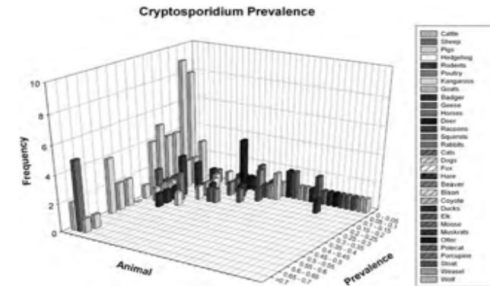
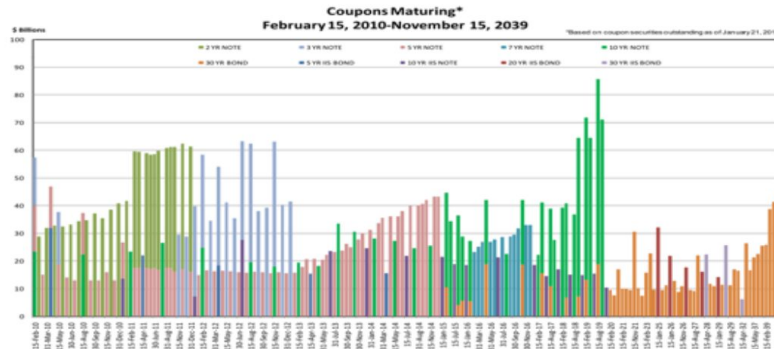
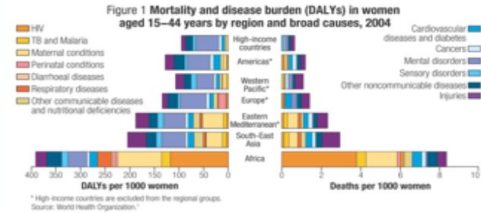
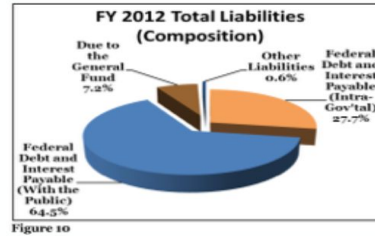
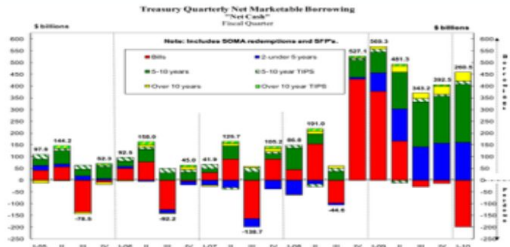


Figure 5.2 Mean prevalence rates of *Cryptosporidium* oocysts by animal species.

Effective Visualizations

- 1. Have graphical integrity**
- 2. Keep it simple**
- 3. Use the right display**
- 4. Use color strategically**

Graphical Integrity

Bonus Assignment

Chart #1:



Issues:

- Tell a Story
 - The vertical axis isn't labeled. We don't know the unit.
- Graphical Integrity
 - The vertical axis does not start from zero
- Graphical Complexity
 - Horizontal and vertical lines are unnecessary (Chartjunk)

Bonus Assignment

Chart #1:



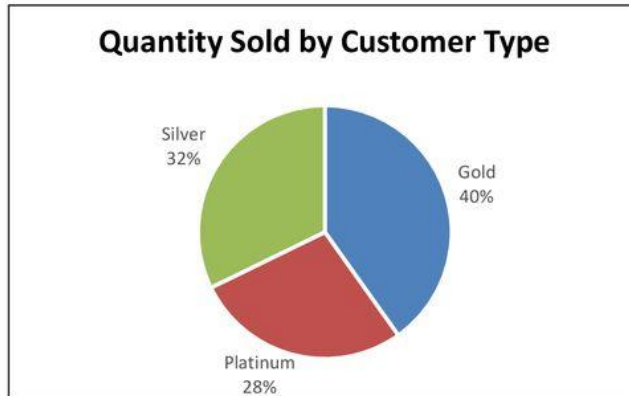
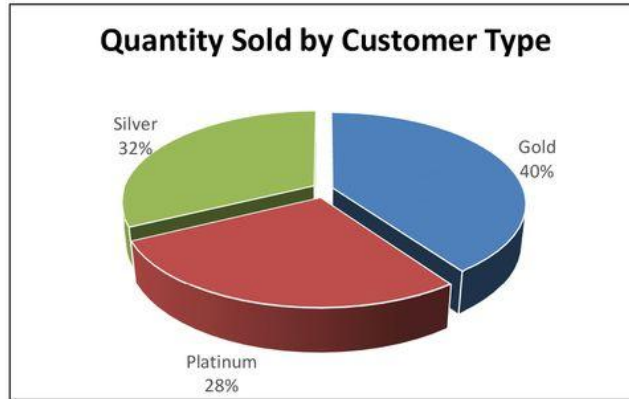
Issues:

- Tell a Story
 - The vertical axis isn't labeled. We don't know the unit.
- Graphical Integrity
 - The vertical axis does not start from zero
- Graphical Complexity
 - Horizontal and vertical lines are unnecessary (Chartjunk)

← This also works

Bonus Assignment

Chart #2:

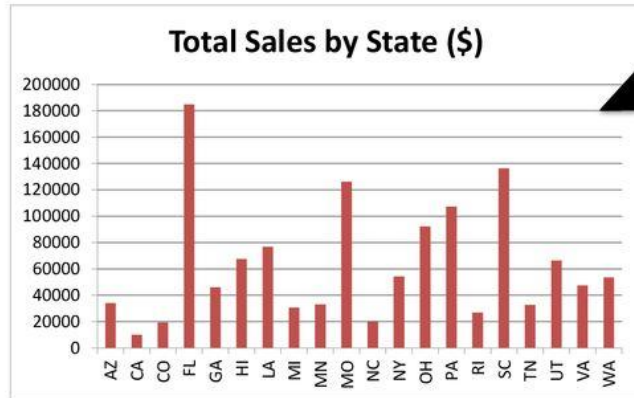
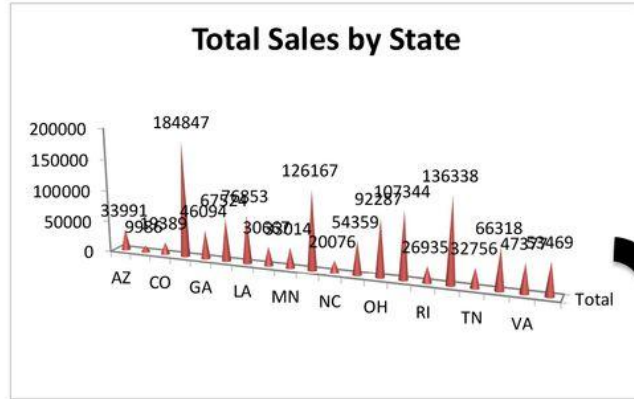


Issues:

- Graphical Integrity
 - The 3D chart makes it difficult to compare the sizes
- Graphical Complexity
 - The 3D chart requires more ink (Chartjunk)

Bonus Assignment

Chart #3:

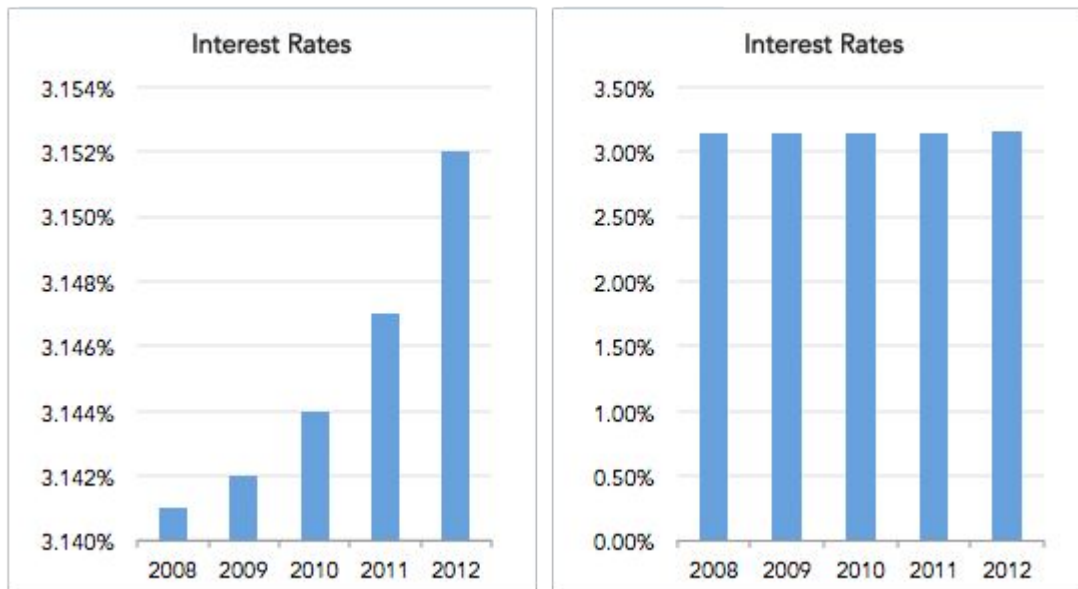


Issues:

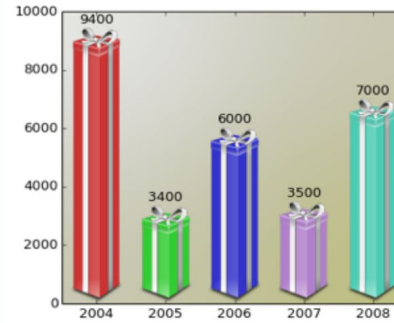
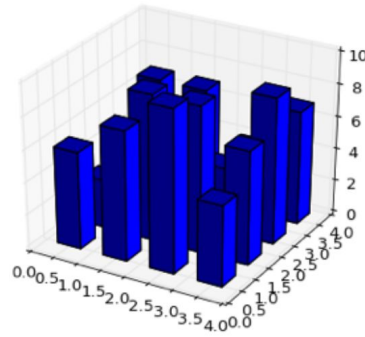
- Tell a Story
 - Vertical axis isn't labeled. We don't know the units
 - Because there are many states to compare, horizontal lines may be helpful
- Graphical Integrity
 - The 3D chart makes it difficult to compare sizes
 - The cone-shaped bars make it even harder to compare sizes
- Graphical Complexity
 - The 3D chart requires more ink (Chartjunk)
 - The number labels are unnecessary

Scale distortion

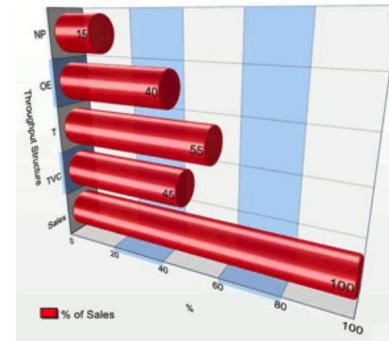
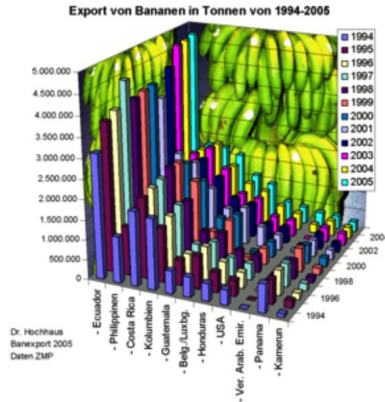
Same Data, Different Y-Axis



Don't!

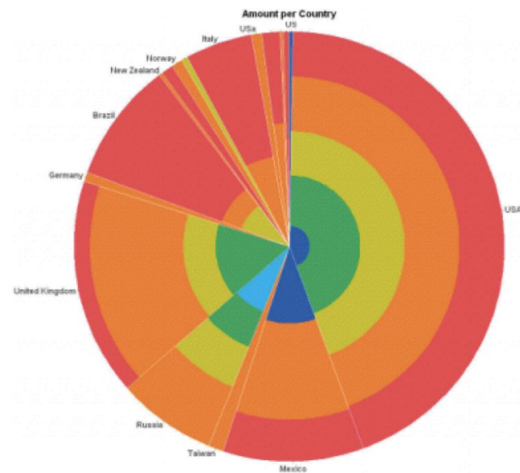
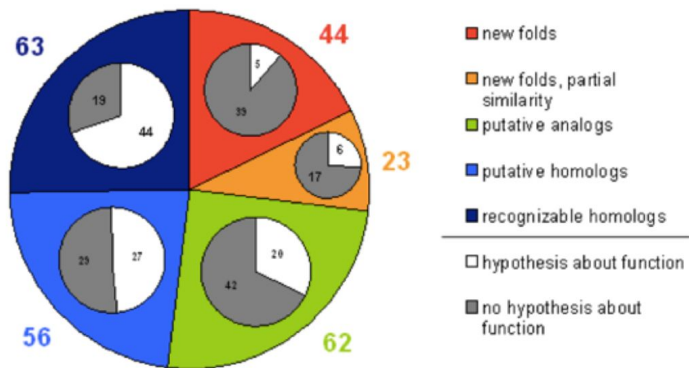


matplotlib gallery

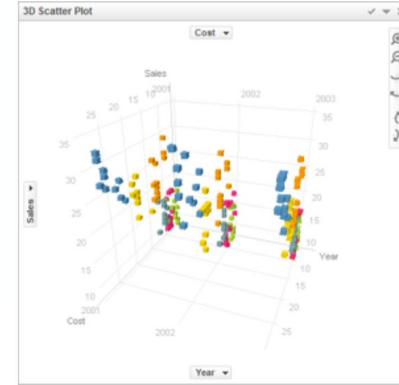
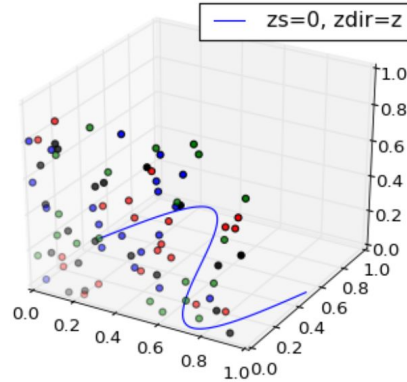
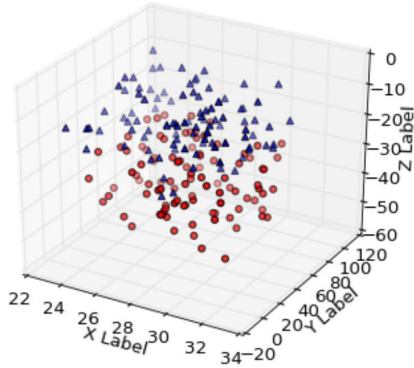


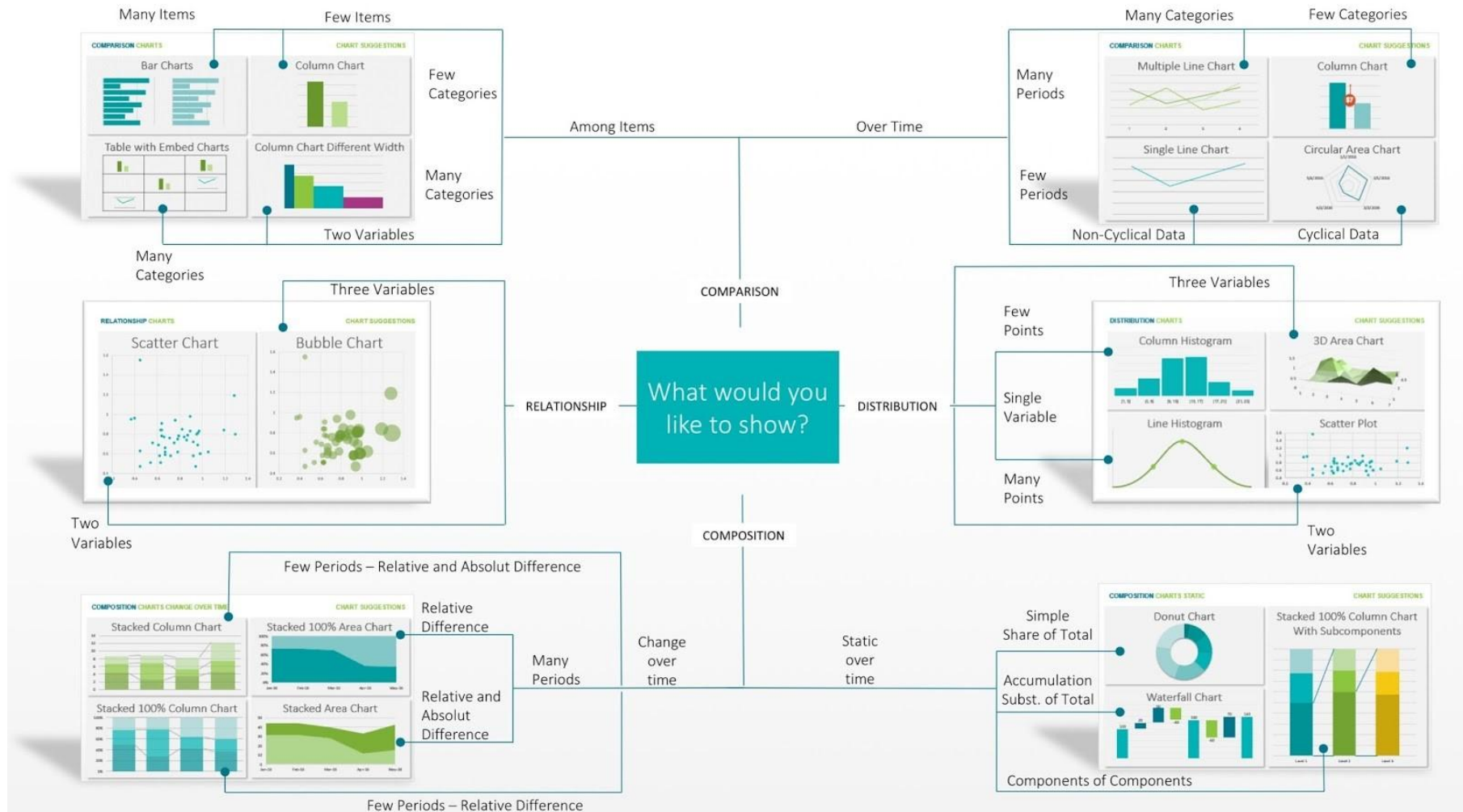
Excel Charts Blog

Don't!



Don't!





Exercise

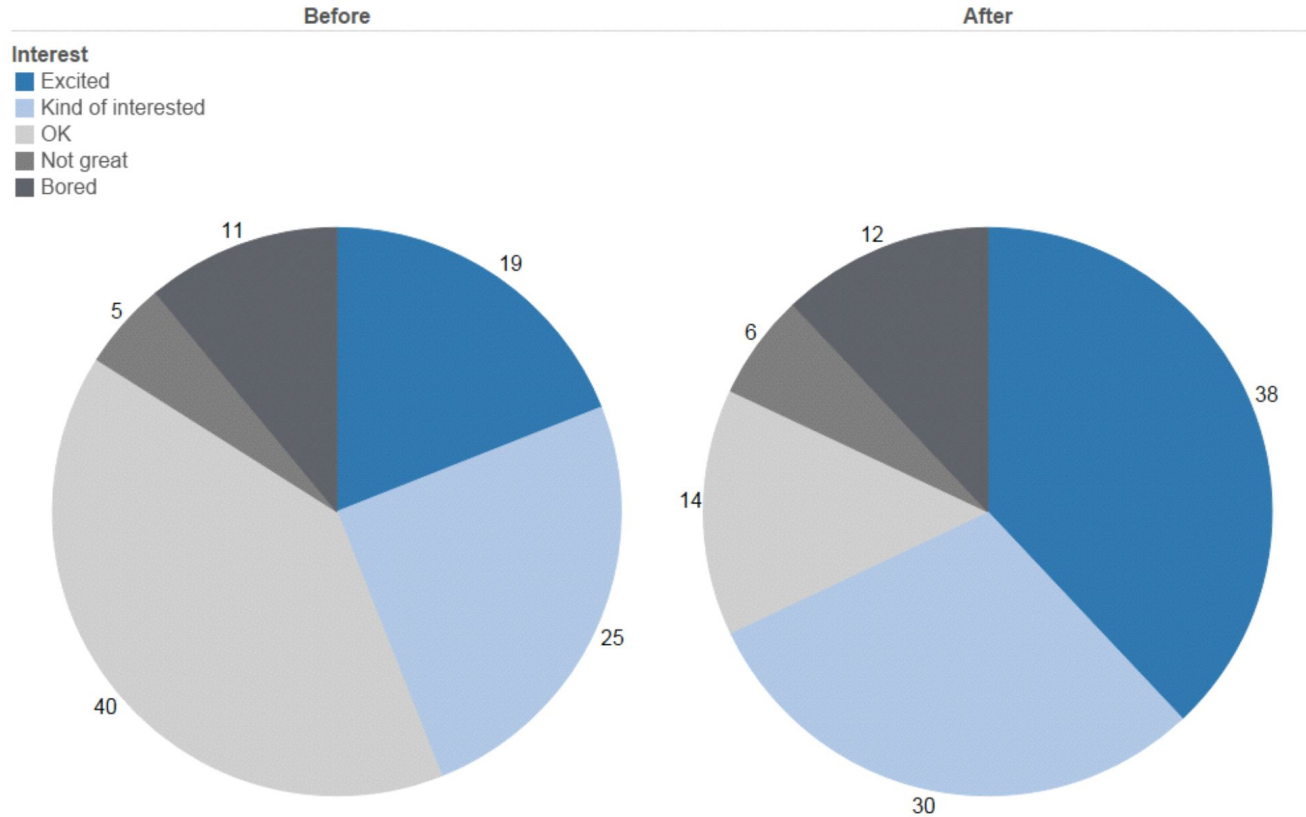
How do you feel about doing science?

Table

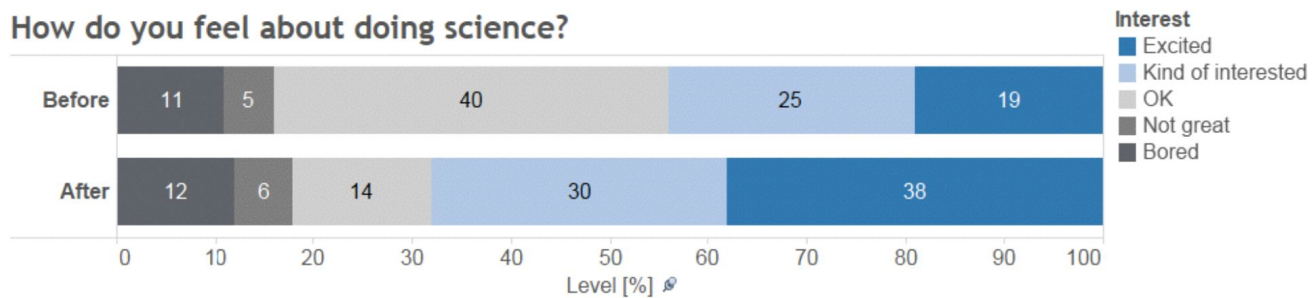
Interest	Before	After
Excited	19	38
Kind of interested	25	30
OK	40	14
Not great	5	6
Bored	11	12

Data courtesy of Cole Nussbaumer

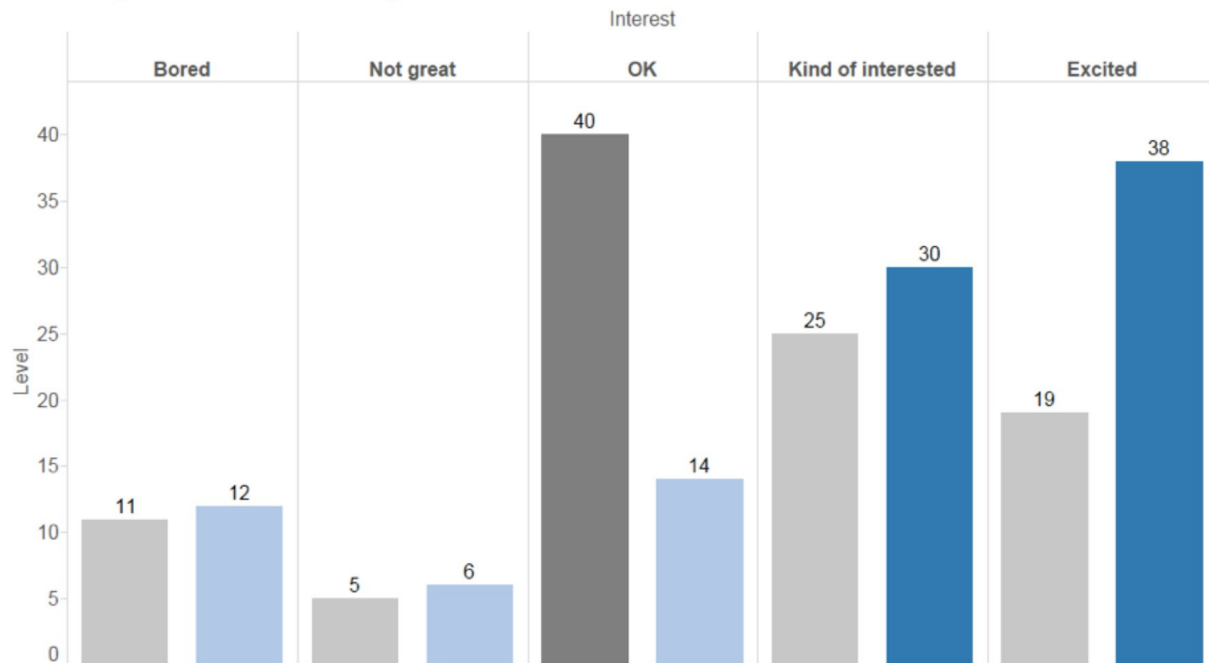
How do you feel about doing science?



How do you feel about doing science?

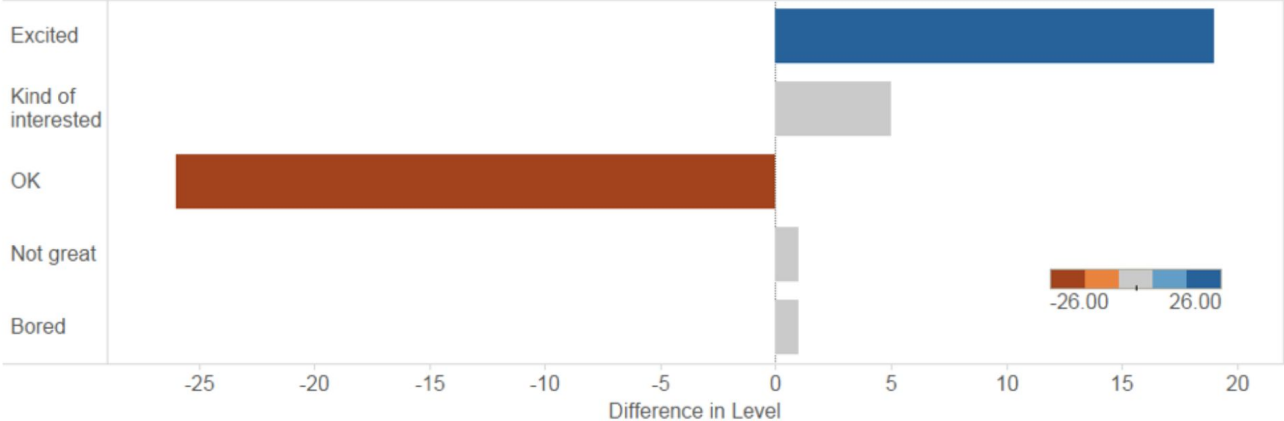


How do you feel about doing science?

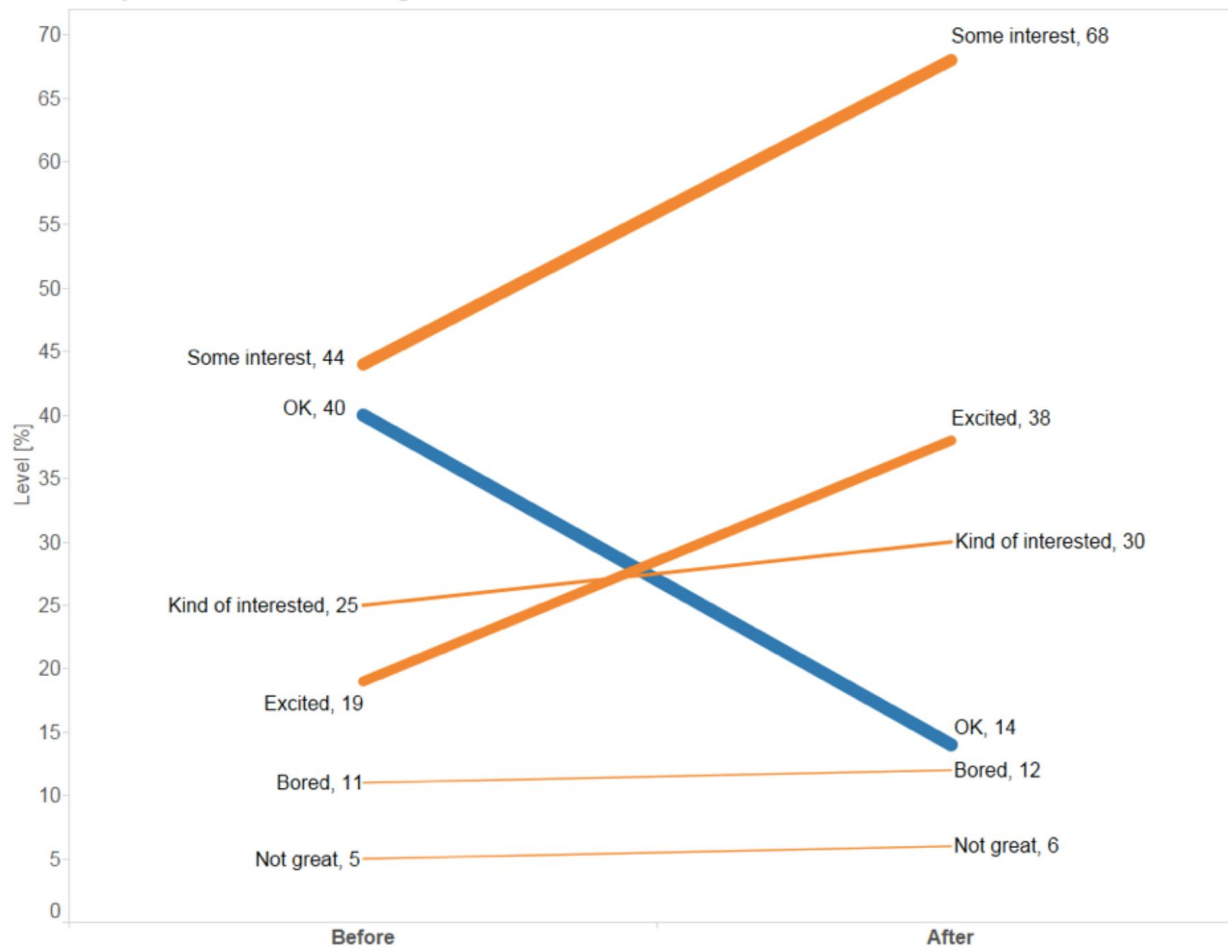


Before the program, the majority of children felt just *OK* about science. After the program, more children were *Kind of interested* and *Excited* about science.

Opinion change to the question: How do you feel about doing science?



How do you feel about doing science?



After the pilot program,

68%

of kids expressed interest towards science,
compared to 44% going into the program.