

Introduction to Data Visualization

...

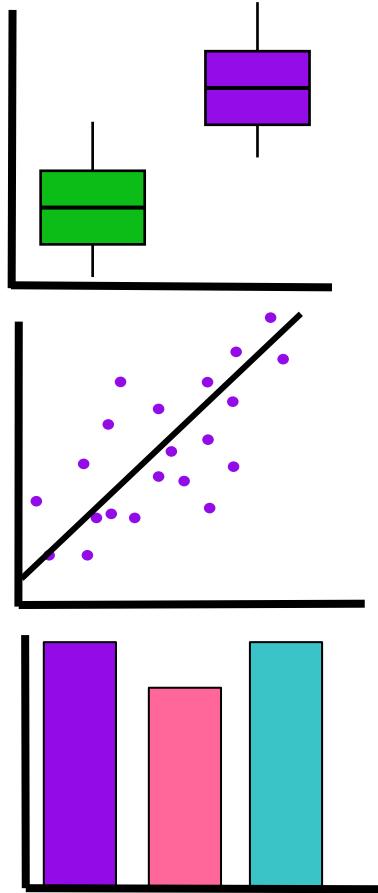
Administrative Items

1. Late submission deadline has passed for assignment 1.
 - a. For questions with grading, email cogs108@gmail.com.
2. Homework 2 was posted over the weekend
 - a. Due date is February 3rd at 11:59
 - b. Section this week will be used for working on this assignment

Data Visualization Types

Graphs

- require at least two scales
- visualize variables and their relationships



Examples: box plots, scatterplots, bar plots

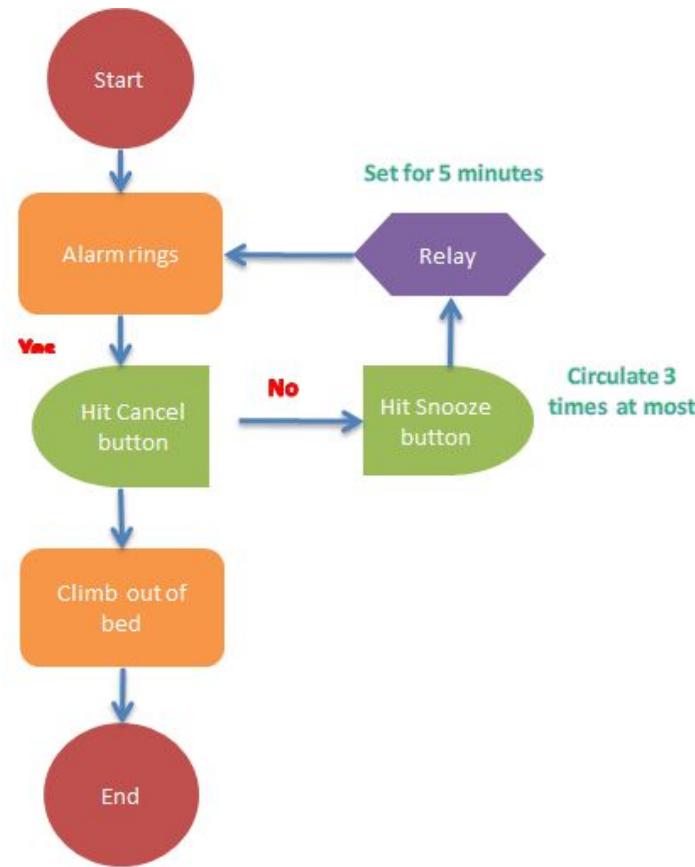
Tables

- display summaries of data
- tabular format

	mean	median
height	68 in	70 in
weight	180 lbs	175 lbs
age	50 yrs	51 yrs

Charts

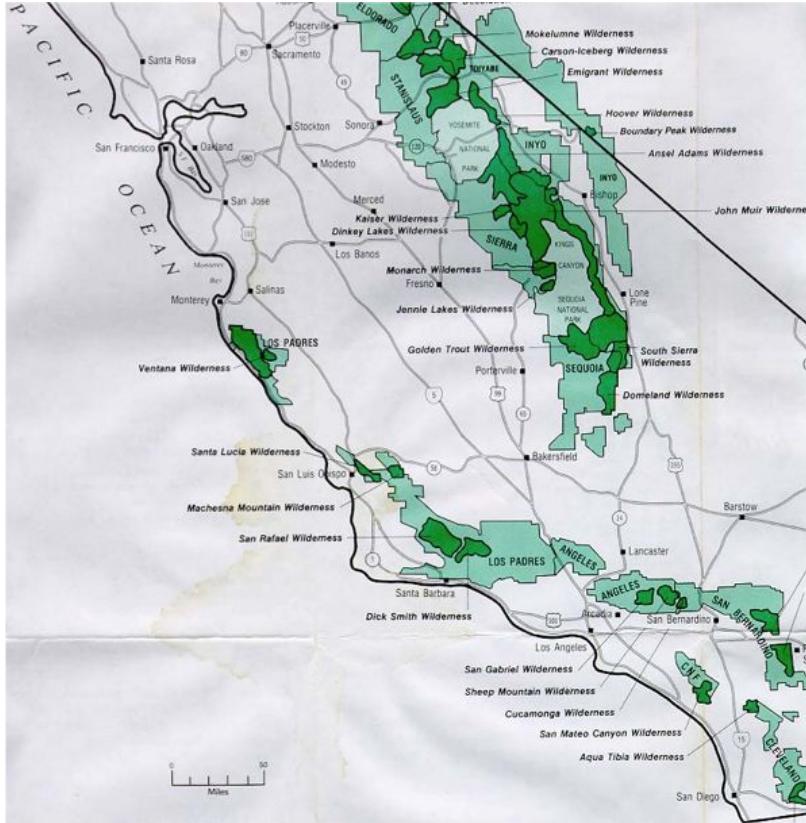
- structure demonstrates relationship
- lines serve as links



Examples: flow charts, family trees, network diagrams

Maps

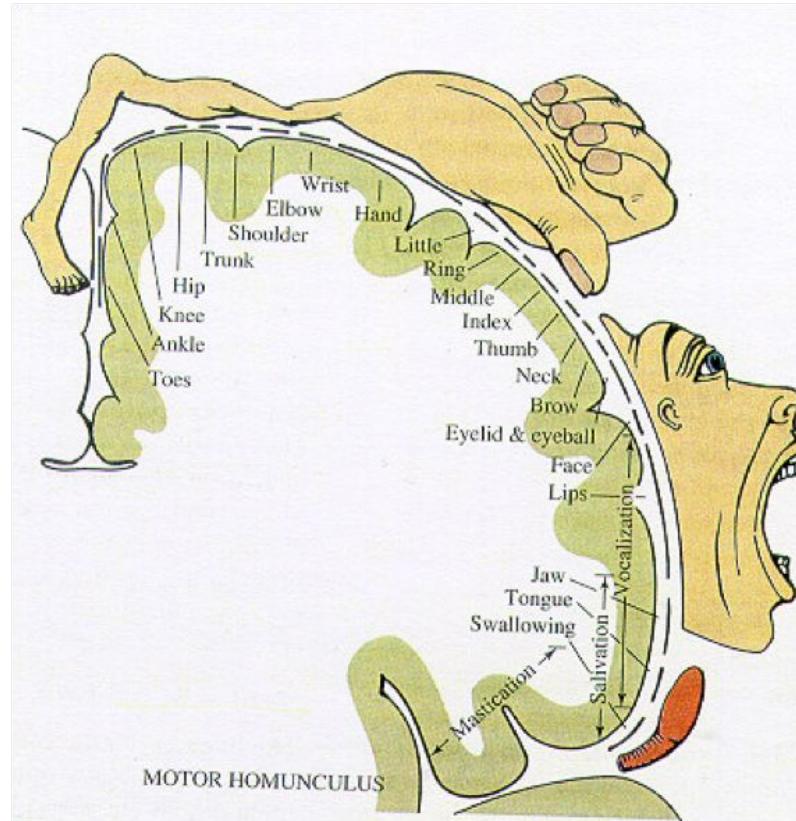
- internal relations determined (in part) by the spatial relations of what is pictured
- labels paired with locations



Examples: physical maps, topographic maps, political maps, maps of census data

Diagrams

- schematic pictures of objects or entities
- parts are *symbolic* (unlike photographs)



Examples: how-to illustrations, figures in a manual

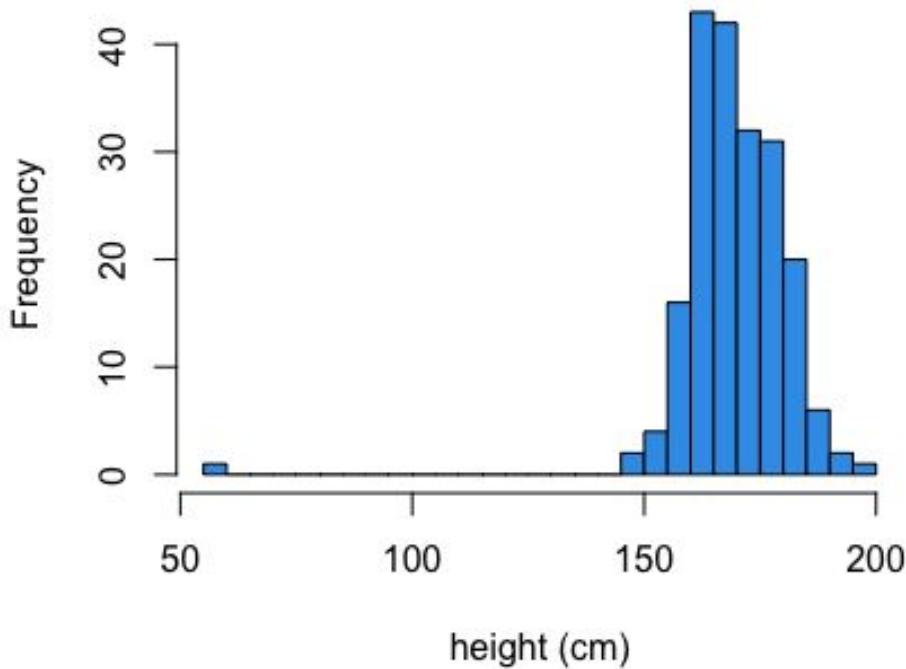
Data Visualization Types: Graphs & Tables

A portion of the
Davis dataset we'll
be working with

	sex	weight	height	repwt	rept
1	M	77	182	77	180
2	F	58	161	51	159
3	F	53	161	54	158
4	M	68	177	70	175
5	F	59	157	59	155
6	M	76	170	76	165
7	M	76	167	77	165
8	M	69	186	73	180
9	M	71	178	71	175
10	M	65	171	64	170
11	M	70	175	75	174
12	F	51	161	52	158
13	F	64	168	64	165
14	F	52	163	57	160
15	F	65	166	66	165
16	M	92	187	101	185
17	F	62	168	62	165
18	M	76	197	75	200
19	F	61	175	61	171
20	M	119	180	124	178
21	F	61	170	61	170
22	M	65	175	66	173
23	M	66	173	70	170
24	F	54	171	59	168

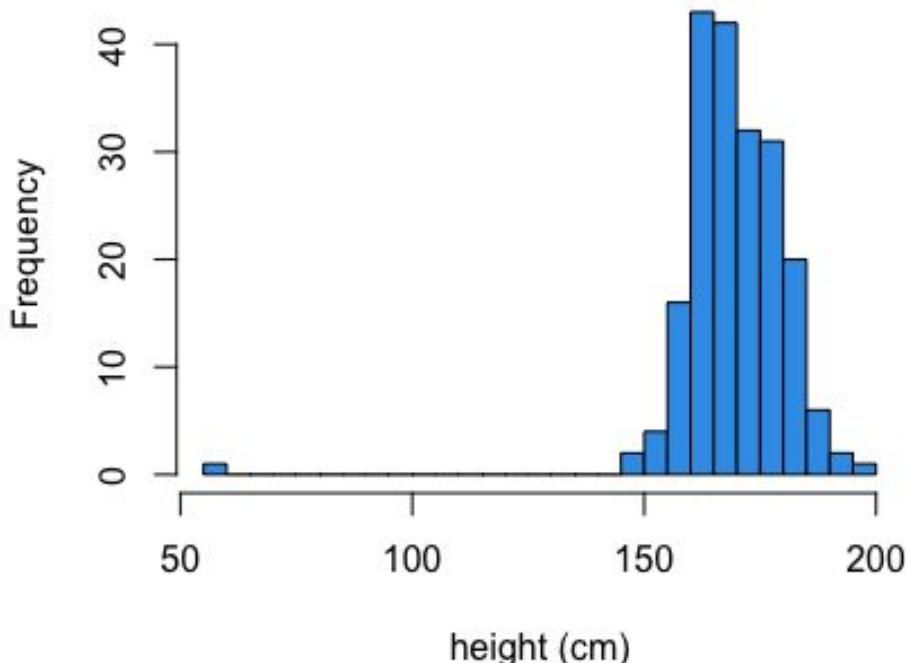
Histograms

Information about
a single set of
numbers



Histograms

Information about
a single set of
numbers

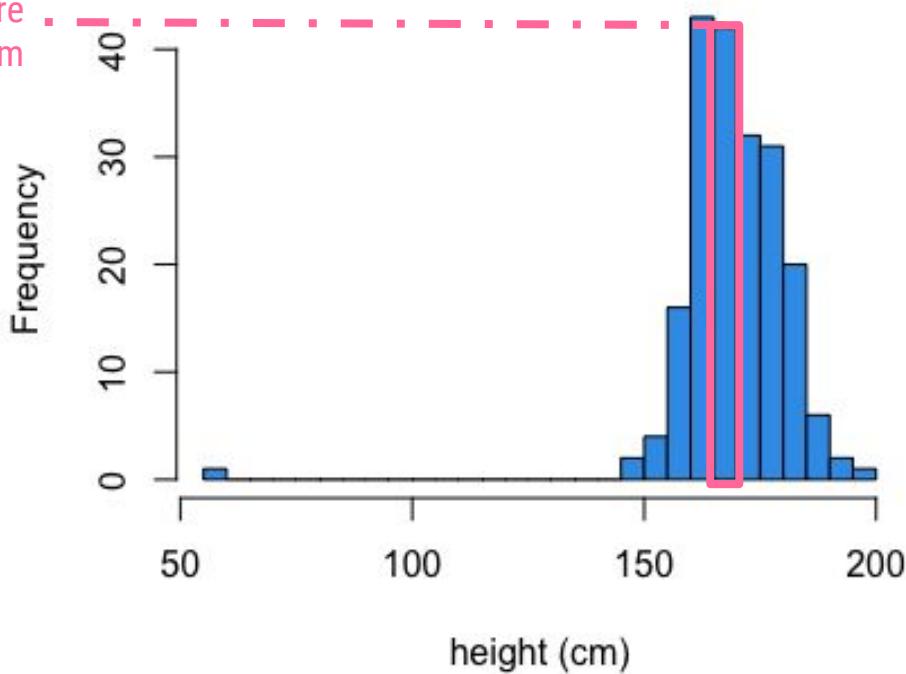


Range of possible height values is easily visualized

Histograms

Information about
a single set of
numbers

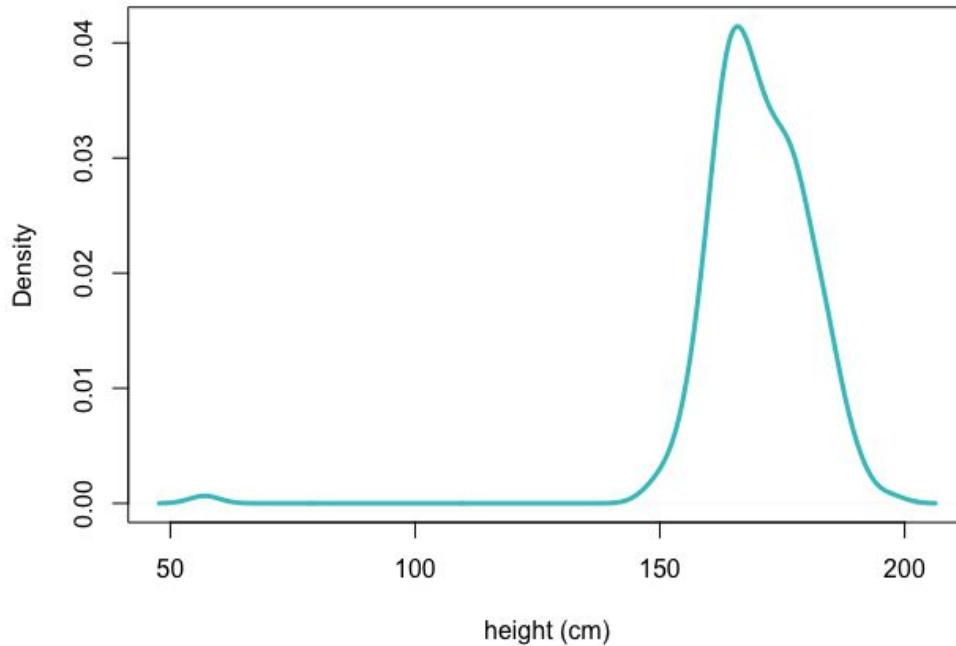
~40 people are
165-170cm



Range of possible height values is easily visualized

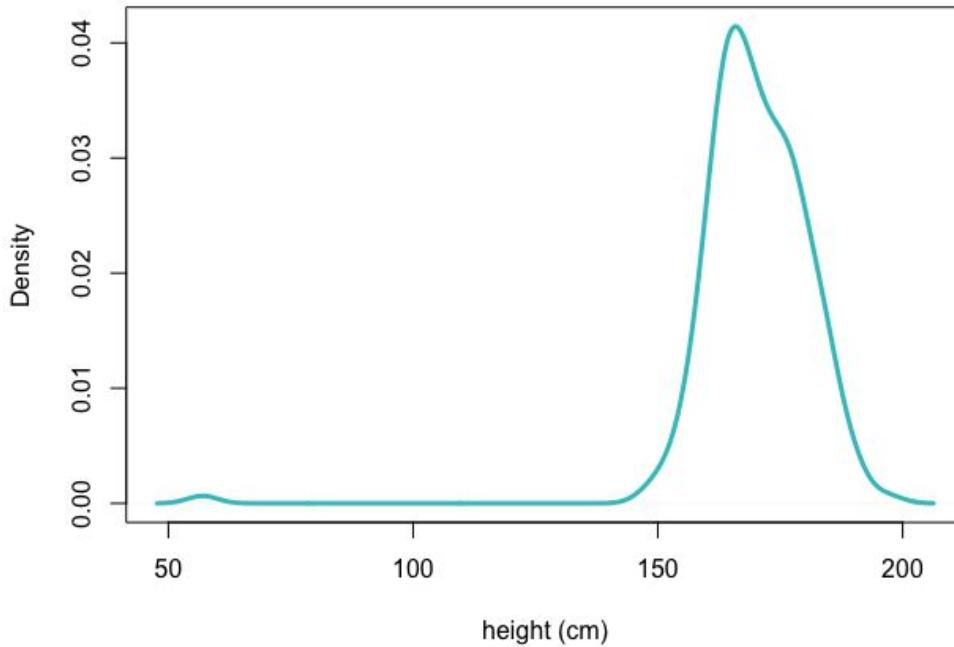
Densityplot

Information about a single set of numbers



Densityplot

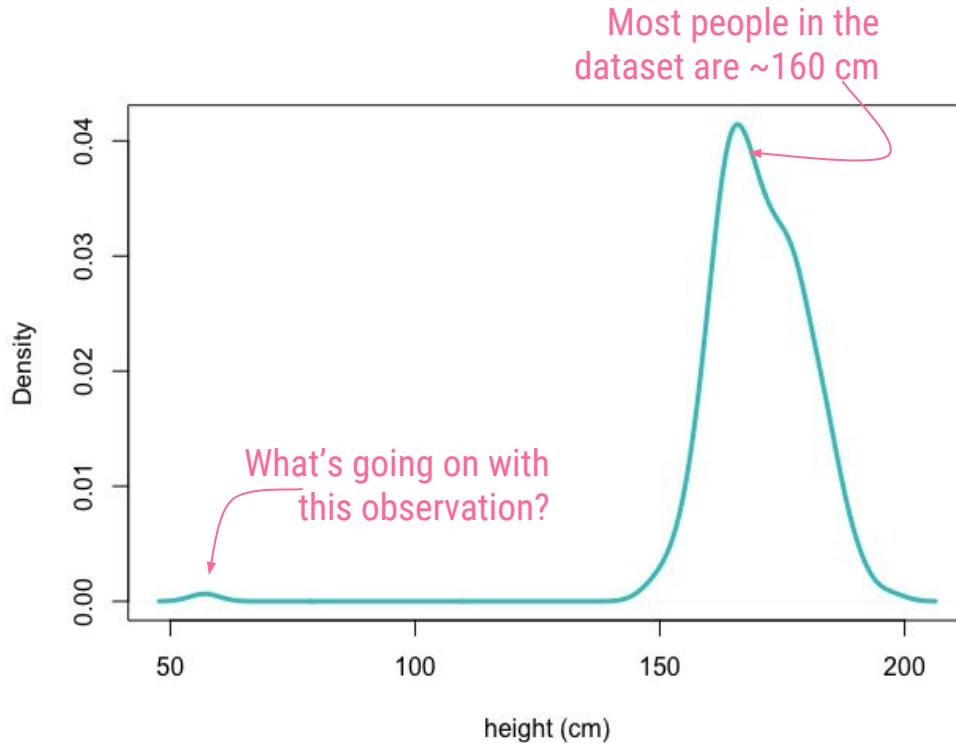
Information about a single set of numbers



Demonstrates the *distribution* of the data
(A smoothed version of a histogram)

Densityplot

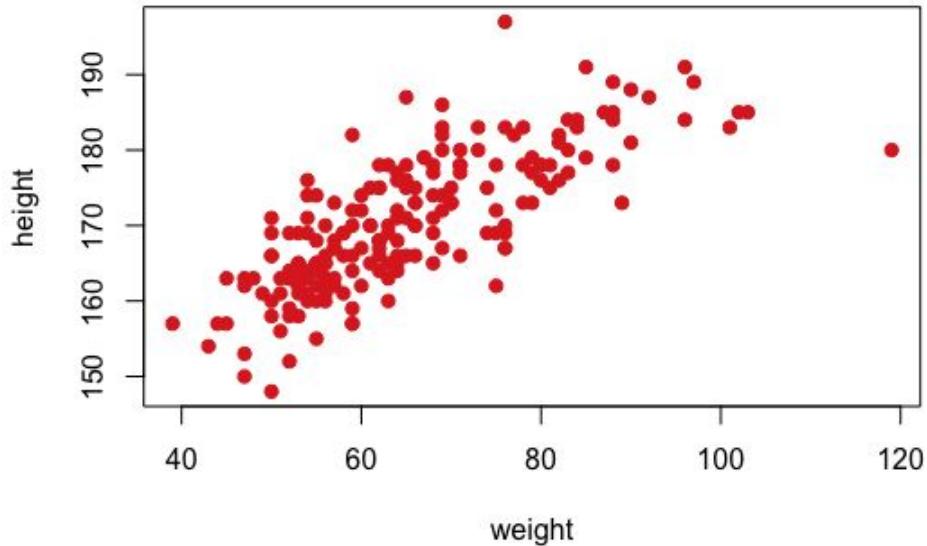
Information about a single set of numbers



Demonstrates the *distribution* of the data
And helps to identify extreme values

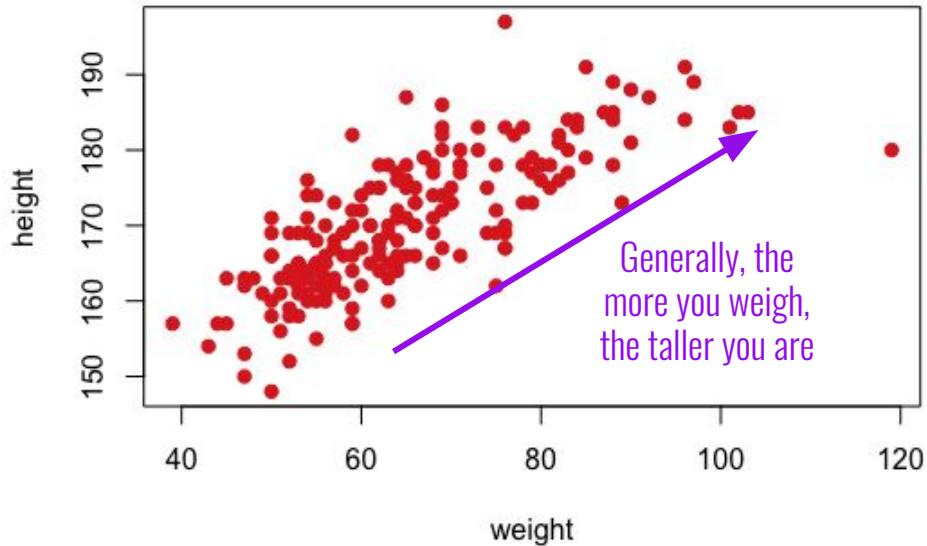
Scatterplot

Relationship between
two numerical
variables



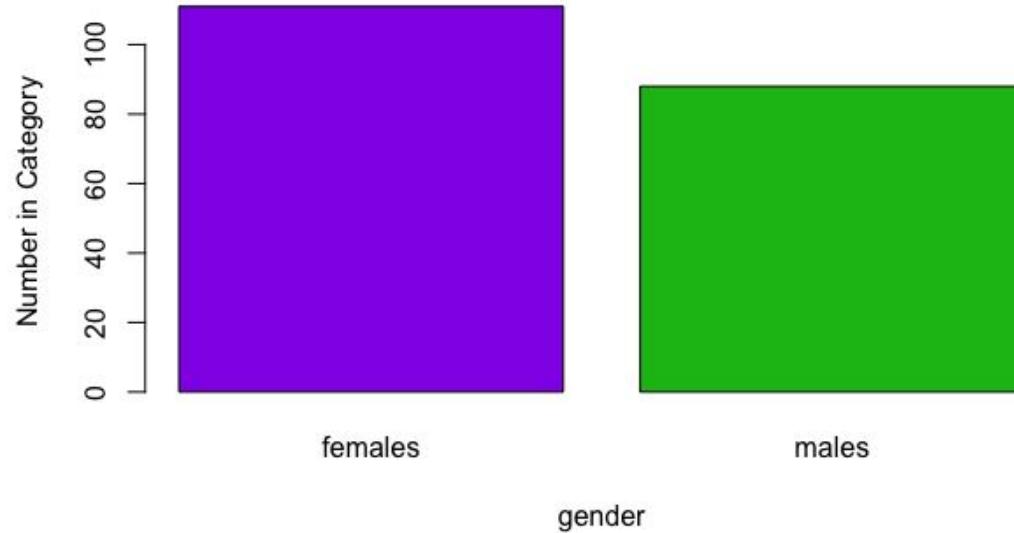
Scatterplot

Relationship between two numerical variables



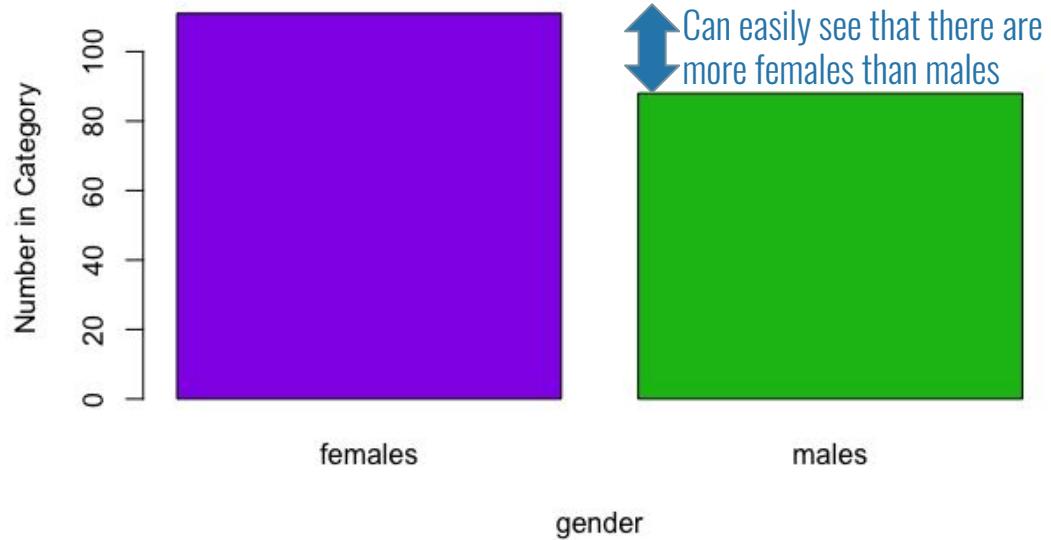
Barplot

Count of levels
within a categorical
variable



Barplot

Count of levels
within a categorical
variable



N=54

5 sets of 10 values were plotted

- 20 graphs (10 pie, 10 bar)

Task: judge which segment or bar is largest

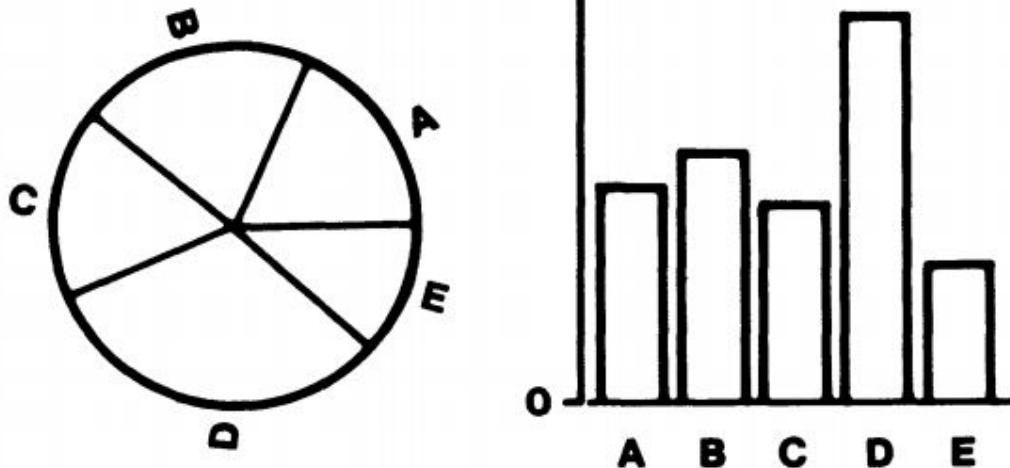
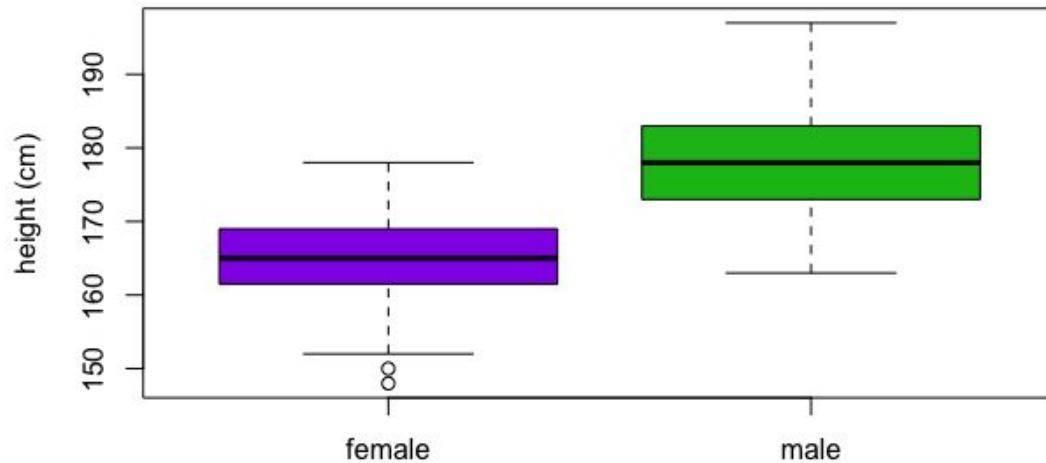


Figure 3. Graphs from position-angle experiment.

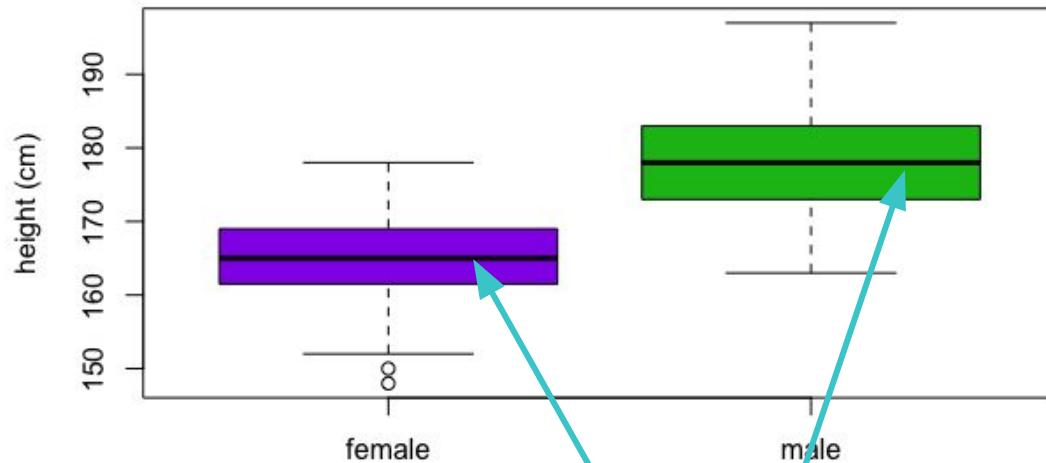
Boxplot

Summary of numerical values across categories



Boxplot

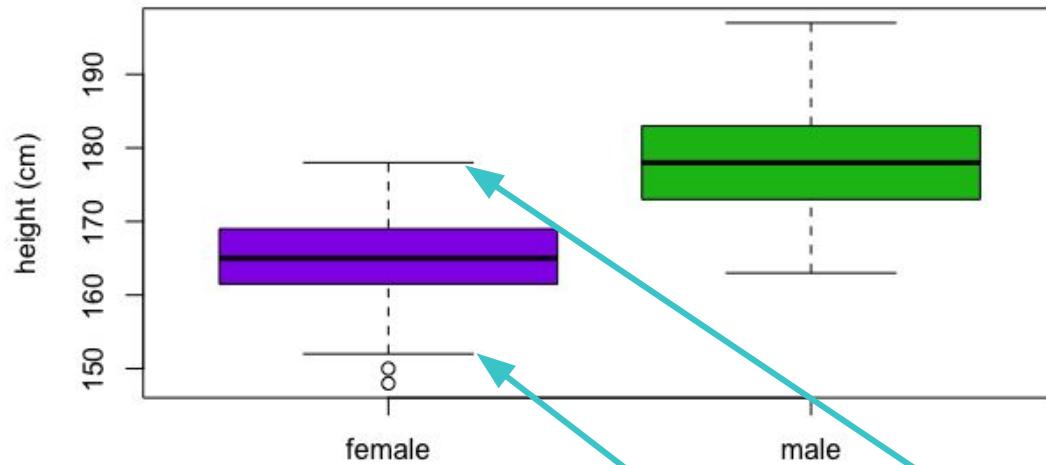
Summary of numerical values across categories



The middle line represents the median & tells you the typical height for females and males

Boxplot

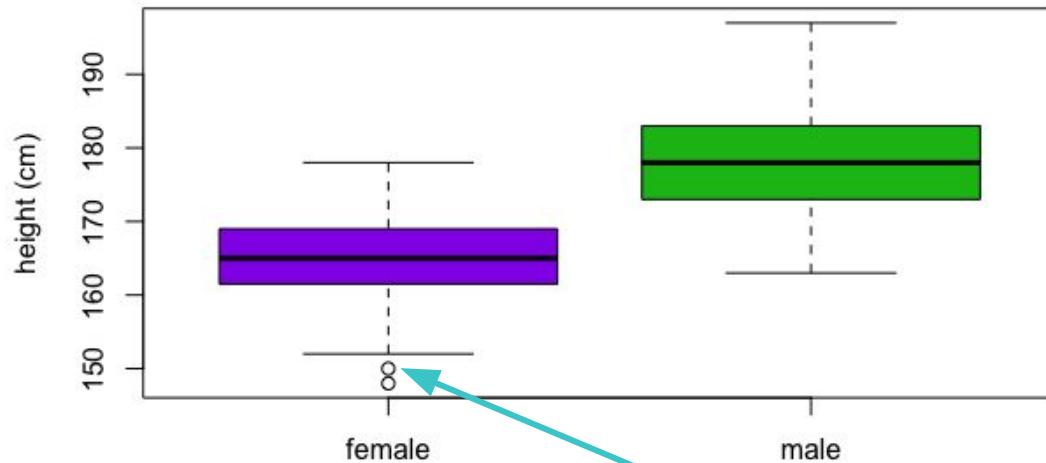
Summary of numerical values across categories



The lines give you an idea of the typical range of values for each category

Boxplot

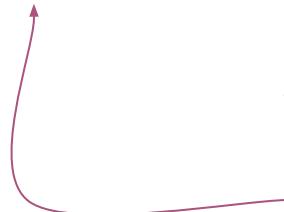
Summary of numerical values across categories



Values outside the typical range are shown as circles. These are known as **outliers**.

Tables

Effective ways to display data summaries



	Number	Percentage (%)
male	111	55.8
female	88	44.2

Table 1: Gender Breakdown Across Davis Dataset

TABLES COUNT AS DATA
VISUALIZATIONS!

Visualization Best Practices

Saliency

Video: Camera 25 (48)
Time: 00:00:00.000 - 00:00:06.000
Participant: User 28
21.75 secs



Extra gentle for the
most sensitive skin.

So gentle on sensitive skin, add the chemicals and moisture
you need to keep you clean if you have diaper rash.

Baby Wipes' unique high-absorbency natural-blend cotton
wipes are 100% cotton soft, extra thick, gel-free protection
for babies with sensitive skin. The chlorine-free materials and
absorbent polymers is non-toxic and anti-irritating. Clinically
tested and pediatrician recommended for babies with allergies
and sensitive skin.

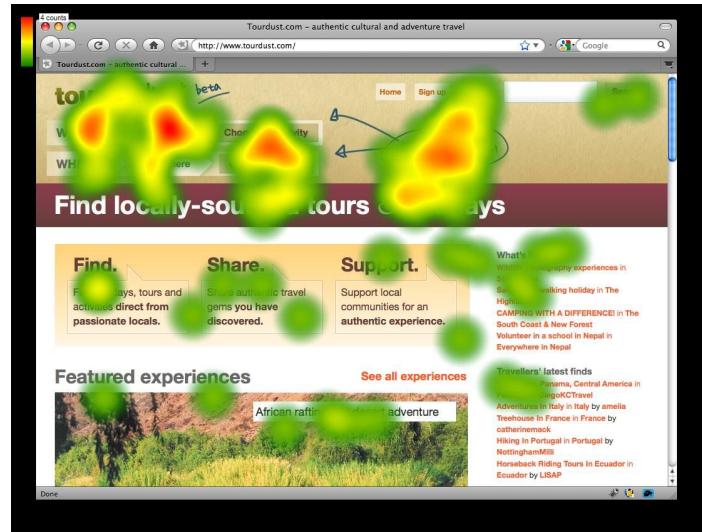
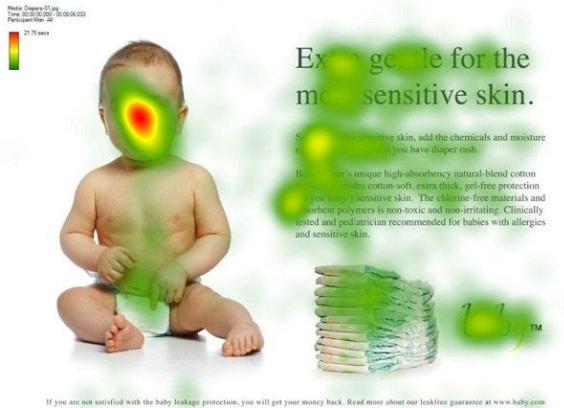


If you are not satisfied with the baby leakage protection, you will get your money back. Read more about our leakfree guarantee at www.baby.com

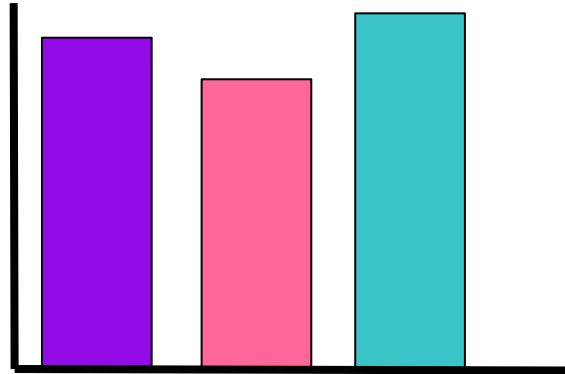
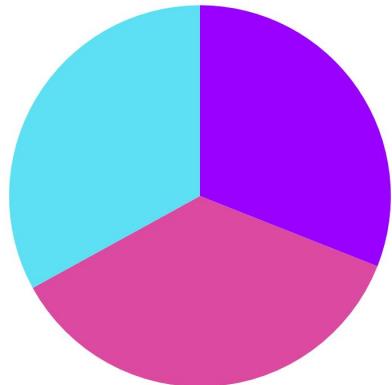
Saliency



Saliency

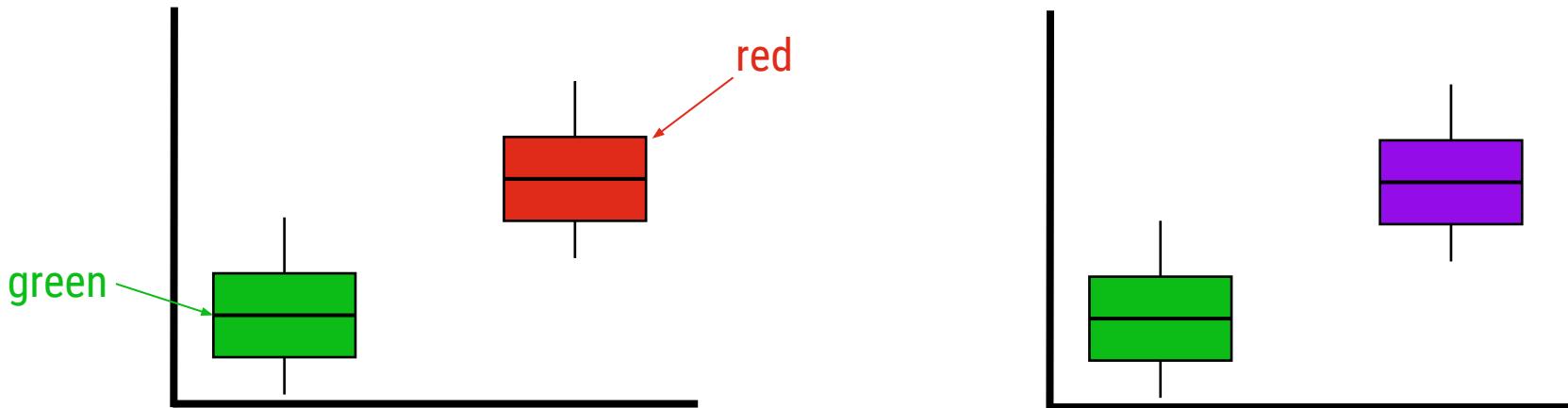


Choose the right type of plot



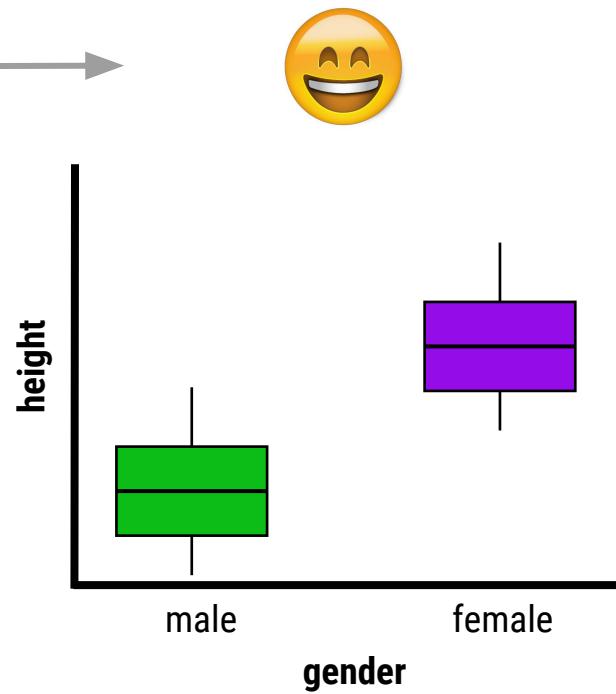
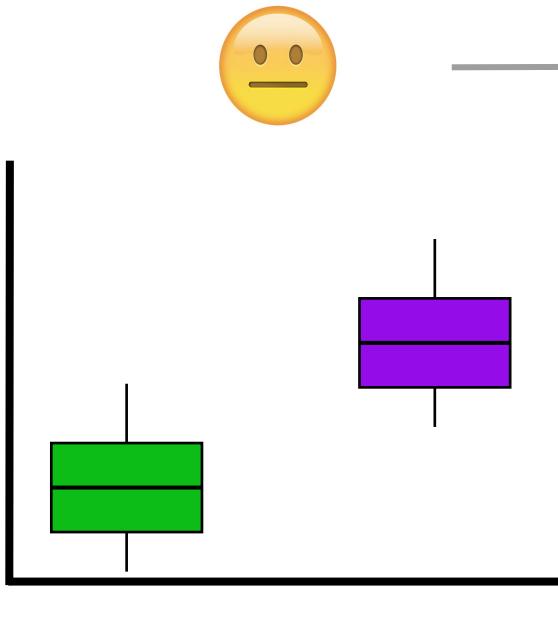
When looking at values, bar charts make it much easier to see the difference between groups!

Be mindful when choosing colors

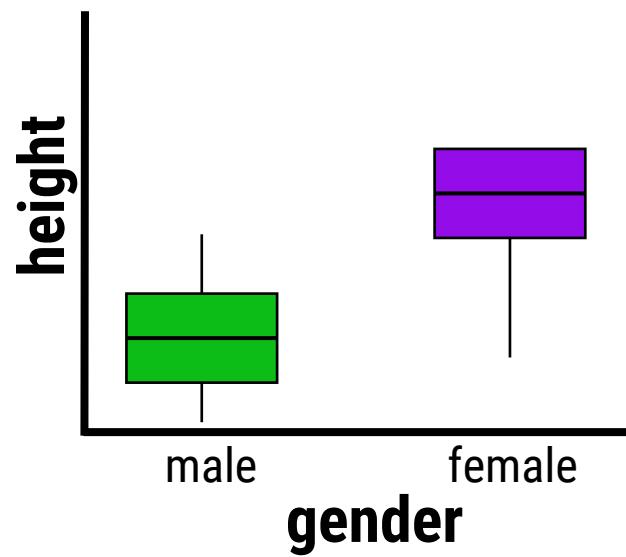
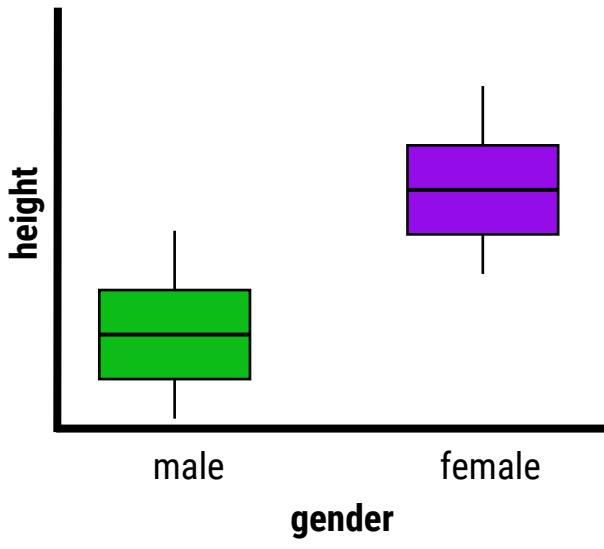


Many color-blind individuals cannot see the difference between red and green.

Label your axes!

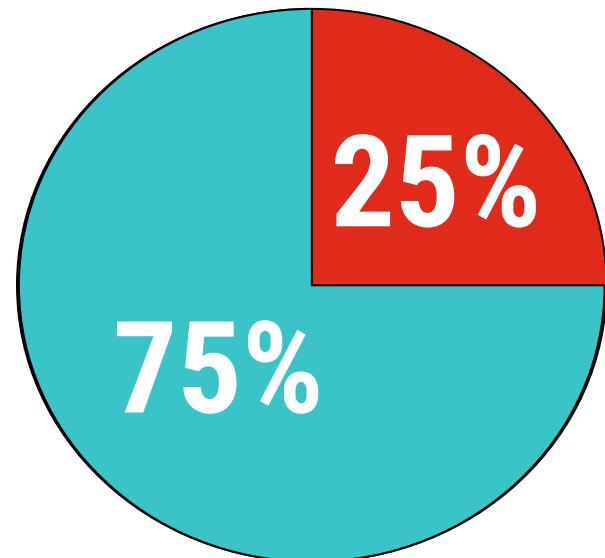
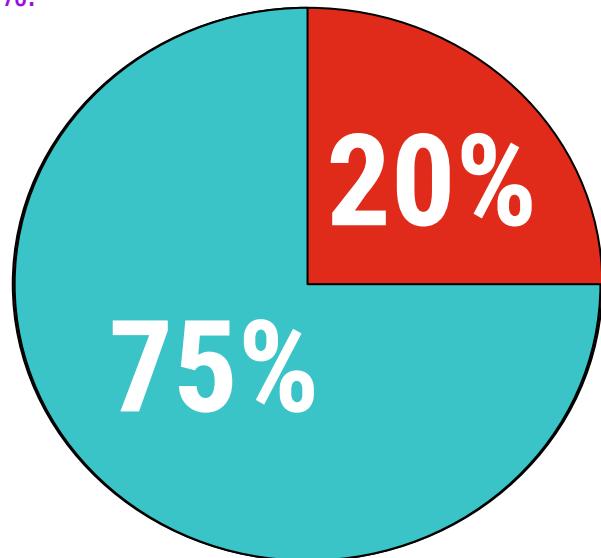


Make sure the text size is big enough!



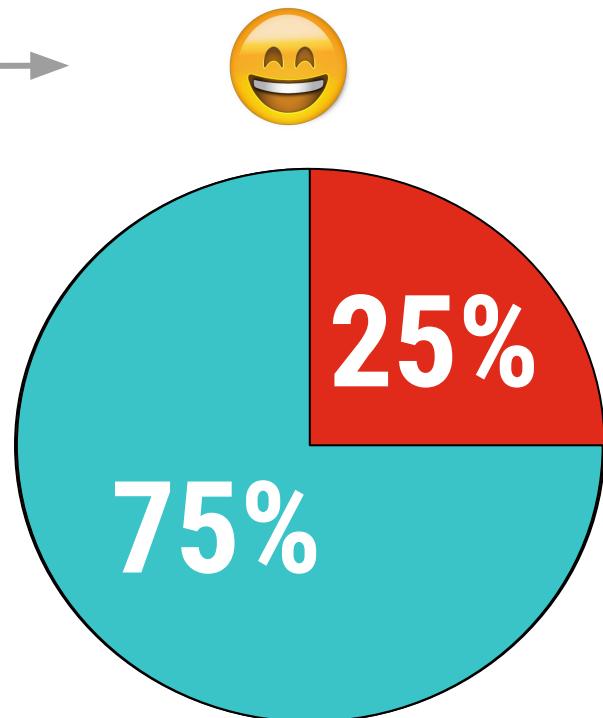
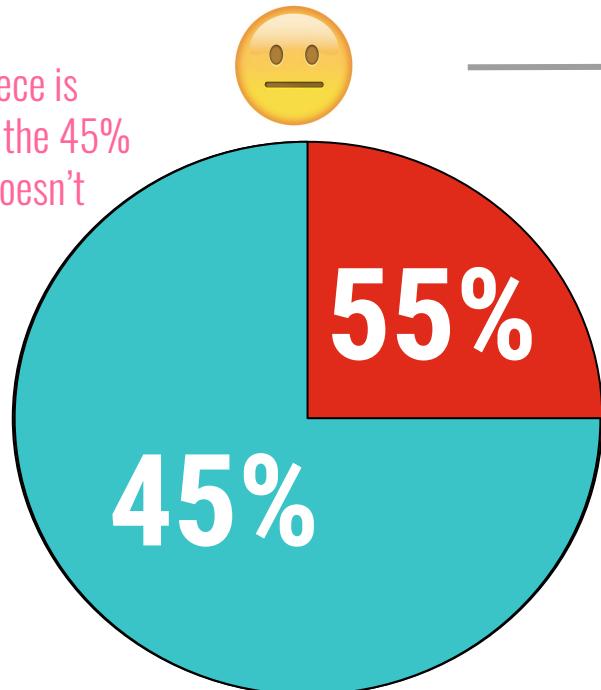
Make sure your numbers add up!

$75 + 20 = 95$...whoops!
Pie charts should always
add up to 100%!



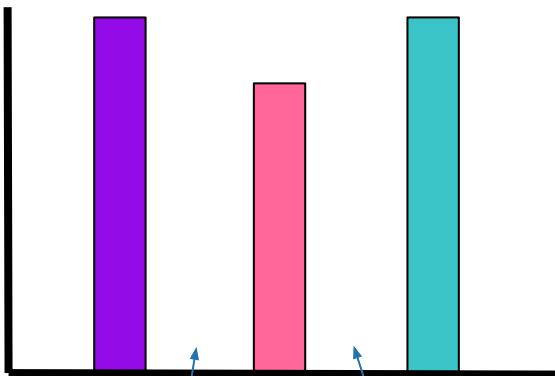
Make sure the numbers and graphic represent the data

That 55% piece is smaller than the 45% piece. That doesn't make sense!

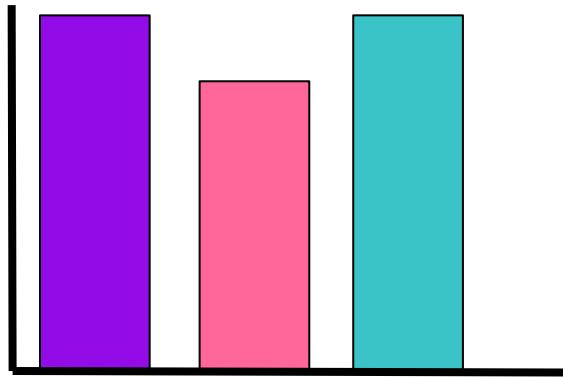


Make comparisons easy on your readers

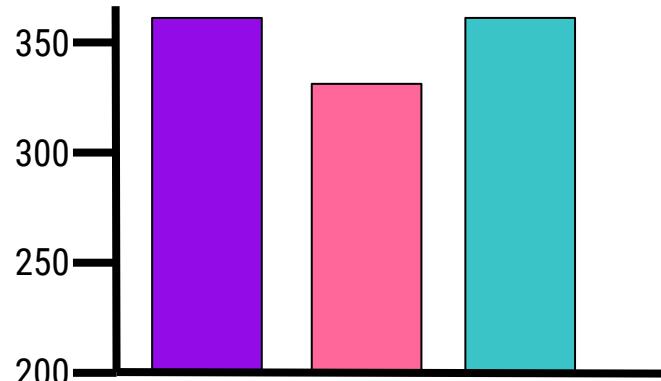
Avoid unnecessary whitespace



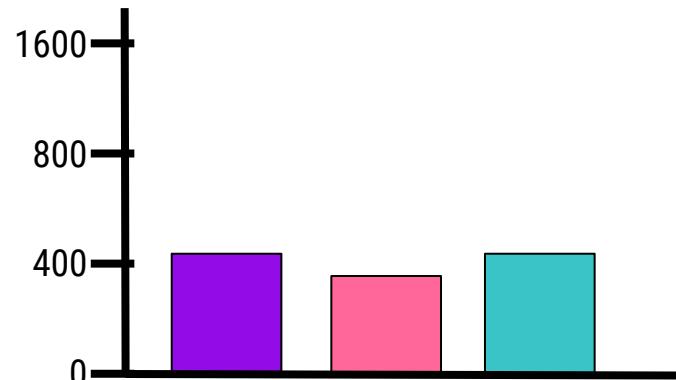
When the bars are spread out,
there is unnecessary whitespace



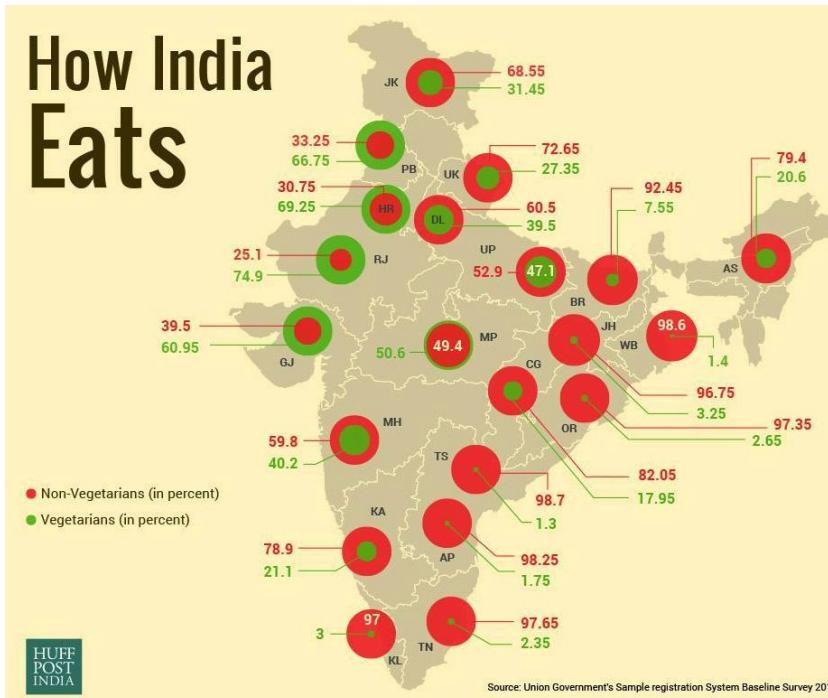
Use y-axes that start at 0 for barplots



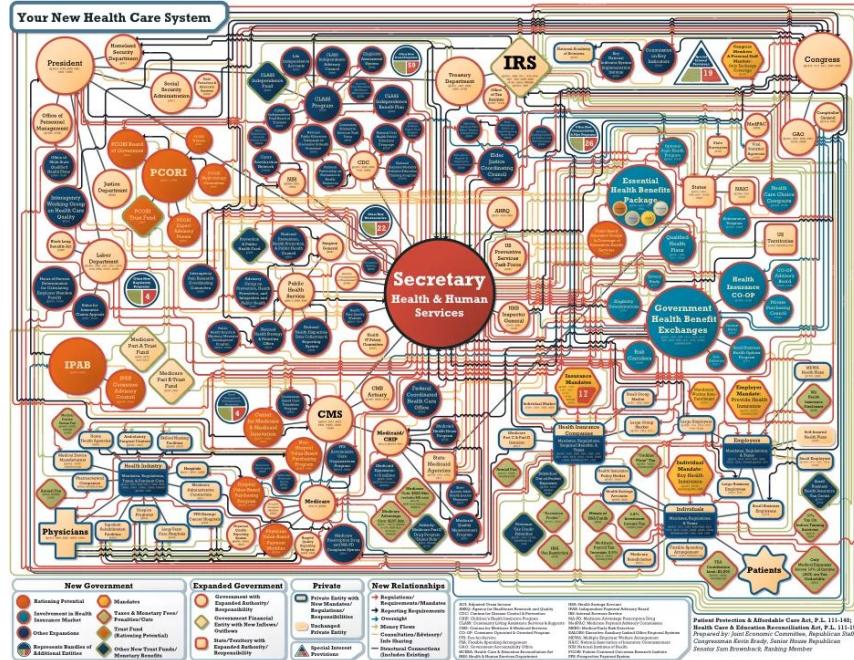
The y-axis starts at 200.
This is misleading & makes differences
look larger than they actually are



Keep it Simple



Keep it Simple



<https://www.jec.senate.gov/public/index.cfm/republicans/committeenews?ID=bb302d88-3d0d-4424-8e33-3c5d2578c2b0>

Allow viewer to make comparison top to bottom



Airport Name	Cleveland Hopkins Intl	William P Hobby	Metropolitan Oakland Intl	San Francisco Intl	Bellingham Intl
Airport Code	CLE	HOU	OAK	SFO	BLI
Mean Arrival Delay	26.15	10.25	10.07	8.86	8.67



Avoid comparisons
across rows

Airport Name	Airport Code	Mean Arrival Delay
	CLE	26.15
William P Hobby	HOU	10.25
Metropolitan Oakland Intl	OAK	10.07
San Francisco Intl	SFO	8.86
Bellingham Intl	BLI	8.67

Compare down columns

Order rows logically



Airport Name	Airport Code	Mean Arrival Delay
Metropolitan Oakland Intl	OAK	10.07
Bellingham Intl	BLI	8.67
	CLE	26.15
San Francisco Intl	SFO	8.86
William P Hobby	HOU	10.25



Airport Name	Airport Code	Mean Arrival Delay
	CLE	26.15
William P Hobby	HOU	10.25
Metropolitan Oakland Intl	OAK	10.07
San Francisco Intl	SFO	8.86
Bellingham Intl	BLI	8.67

longest delay

shortest delay

Order columns logically



A diagram showing a transformation from a state of confusion (represented by a confused emoji above a table) to a state of clarity (represented by a smiling emoji above a sorted table). A horizontal arrow points from left to right, indicating the process of ordering columns logically.

Mean Arrival Delay	Airport Code	Airport Name
26.15	CLE	
10.25	HOU	William P Hobby
10.07	OAK	Metropolitan Oakland Intl
8.86	SFO	San Francisco Intl
8.67	BLI	Bellingham Intl



A diagram showing a transformation from a state of confusion (represented by a confused emoji above a table) to a state of clarity (represented by a smiling emoji above a sorted table). A horizontal arrow points from left to right, indicating the process of ordering columns logically.

Airport Name	Airport Code	Mean Arrival Delay
	CLE	26.15
William P Hobby	HOU	10.25
Metropolitan Oakland Intl	OAK	10.07
San Francisco Intl	SFO	8.86
Bellingham Intl	BLI	8.67

Limit the number of rows and columns



	Warren	Median	UrbanPop	Rape
Alabama	13.2	256	58	21.2
Alaska	10.0	263	48	44.2
Arizona	6.1	241	60	37.0
Arkansas	8.6	190	50	19.5
California	9.0	276	91	40.5
Colorado	7.9	204	78	38.7
Connecticut	3.3	110	77	11.1
Delaware	5.9	238	72	15.8
Florida	15.4	305	80	31.9
Georgia	17.4	231	60	35.8
Hawaii	5.3	48	60	20.3
Idaho	2.6	120	54	14.2
Illinois	10.4	249	83	24.0
Indiana	7.2	113	65	21.0
Iowa	2.2	56	57	11.3
Kansas	6.0	115	66	18.0
Kentucky	4.7	160	52	18.3
Louisiana	15.4	249	66	22.3
Maine	2.1	83	91	7.8
Maryland	11.3	300	67	27.8
Massachusetts	4.4	148	95	16.3
Michigan	12.1	255	74	35.1
Minnesota	2.7	72	65	14.8
Mississippi	9.1	259	44	17.7
Missouri	8.0	178	70	28.2
Montana	6.0	109	53	16.4
Nebraska	4.3	102	62	16.5
Nevada	12.2	252	81	46.0
New Hampshire	2.1	57	56	9.5
New Jersey	7.4	159	89	18.8
New Mexico	11.4	260	70	25.1
New York	11.1	254	96	26.1
North Carolina	13.0	337	45	16.1
North Dakota	0.8	45	44	7.3
Ohio	7.3	120	75	21.4
Oklahoma	6.6	151	68	20.0
Oregon	4.9	159	97	29.9
Pennsylvania	6.3	169	73	14.8
Rhode Island	3.4	174	87	8.3
South Carolina	14.4	279	48	22.5
South Dakota	3.8	86	40	12.8
Tennessee	13.2	188	69	26.9
Texas	12.7	201	80	25.6
Utah	5.2	167	60	25.9
Vermont	2.2	48	33	11.2
Virginia	8.5	196	63	20.7
Washington	4.0	145	73	26.2
West Virginia	5.7	81	59	9.3
Wisconsin	2.6	53	66	10.8
Wyoming	6.8	161	60	15.6



Too many rows!



Airport Name

Airport Name	Airport Code	Mean Arrival Delay
	CLE	26.15
William P Hobby	HOU	10.25
Metropolitan Oakland Intl	OAK	10.07
San Francisco Intl	SFO	8.86
Bellingham Intl	BLI	8.67

Include informative labels



 Needs better labels!

AN	AC	MAD
	CLE	26.15
William P Hobby	HOU	10.25
Metropolitan Oakland Intl	OAK	10.07
San Francisco Intl	SFO	8.86
Bellingham Intl	BLI	8.67

Airport Name	Airport Code	Mean Arrival Delay
	CLE	26.15
William P Hobby	HOU	10.25
Metropolitan Oakland Intl	OAK	10.07
San Francisco Intl	SFO	8.86
Bellingham Intl	BLI	8.67

Be mindful of significant digits



Airport Name	Airport Code	Mean Arrival Delay
	CLE	26.15145
William P Hobby	HOU	10.258923
Metropolitan Oakland Intl	OAK	10.07987
San Francisco Intl	SFO	8.8688173
Bellingham Intl	BLI	8.671234

Too many digits

Airport Name	Airport Code	Mean Arrival Delay
	CLE	26.15
William P Hobby	HOU	10.25
Metropolitan Oakland Intl	OAK	10.07
San Francisco Intl	SFO	8.86
Bellingham Intl	BLI	8.67

Format table so it can be quickly understood



Table 1: Airports with the longest delay in arrival time for planes leaving the Pacific Northwest in 2014

Airport Name	Airport Code	Mean Arrival Delay
	CLE	26.15
William P Hobby	HOU	10.25
Metropolitan Oakland Intl	OAK	10.07
San Francisco Intl	SFO	8.86
Bellingham Intl	BLI	8.67

Include a good caption



Table 1: Airports with the longest delay in arrival time for planes leaving the Pacific Northwest in 2014

Airport Name	Airport Code	Mean Arrival Delay
	CLE	26.15
William P Hobby	HOU	10.25
Metropolitan Oakland Intl	OAK	10.07
San Francisco Intl	SFO	8.86
Bellingham Intl	BLI	8.67

Include the source of the data



Table 1: Airports with the longest delay in arrival time for planes leaving the Pacific Northwest in 2014

Airport Name	Airport Code	Mean Arrival Delay
	CLE	26.15
William P Hobby	HOU	10.25
Metropolitan Oakland Intl	OAK	10.07
San Francisco Intl	SFO	8.86
Bellingham Intl	BLI	8.67

Source: Chester Ismay's `pnwflights14` R package

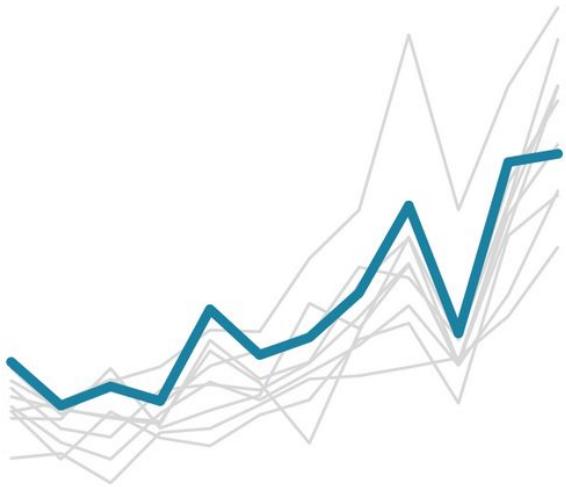
Color Considerations

Gray is your friend.

NOT IDEAL



BETTER



Let the plot highlight values. Use colors for categories.



PEOPLE IN GROUP A



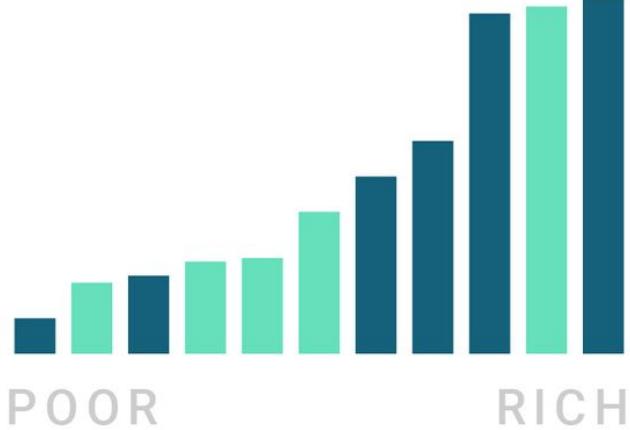
PEOPLE IN GROUP B



NOT IDEAL

BETTER

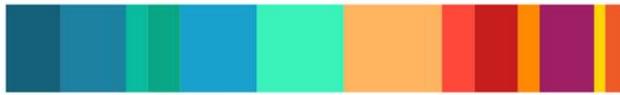
PEOPLE IN GROUP
A B



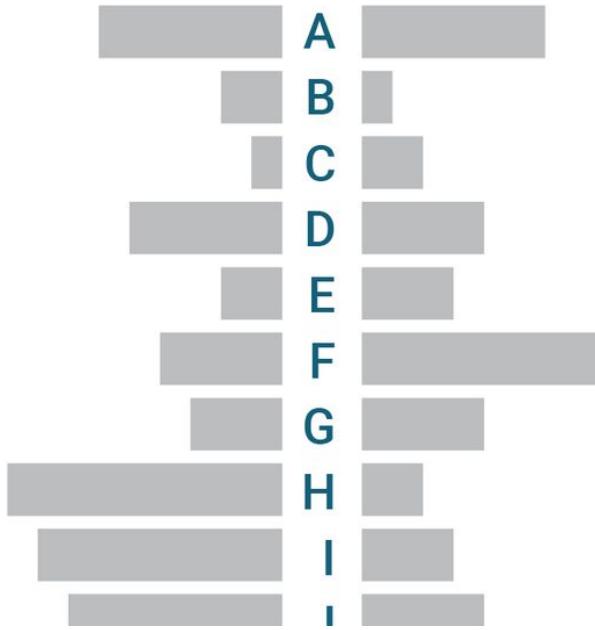
Gradients are great for patterns (i.e. population on a map), but aren't effective at highlighting your point. Consider showing your most important values with bars, position (like in a dot plot) or even areas, and to use colors to only show categories. Consider using

More than 7 colors suggests you need a different chart.

NOT IDEAL

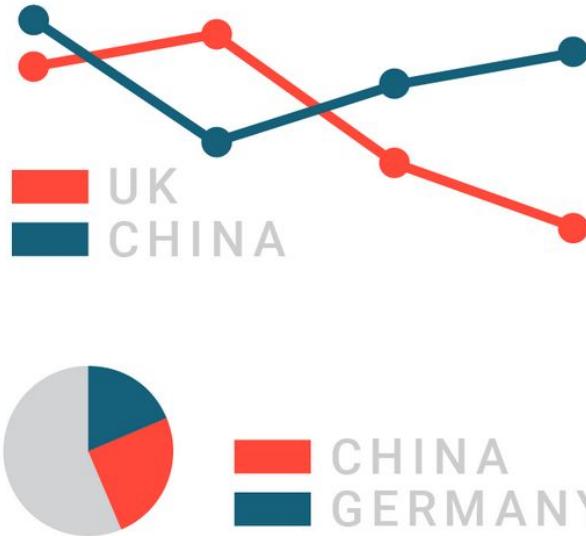


BETTER



Color consistency helps viewers.

NOT IDEAL

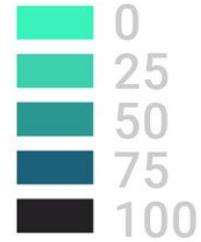


BETTER



Always ensure that viewers know what the colors represent.

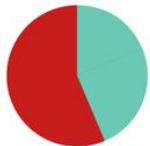
COLOR KEY



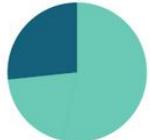
SHARE OF
PEOPLE IN
CHINA AND
GERMANY

Choose intuitive colors

NOT IDEAL



GOOD
BAD

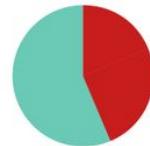


FOREST
LAKE

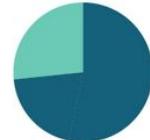


FEMALE
MALE

BETTER



GOOD
BAD



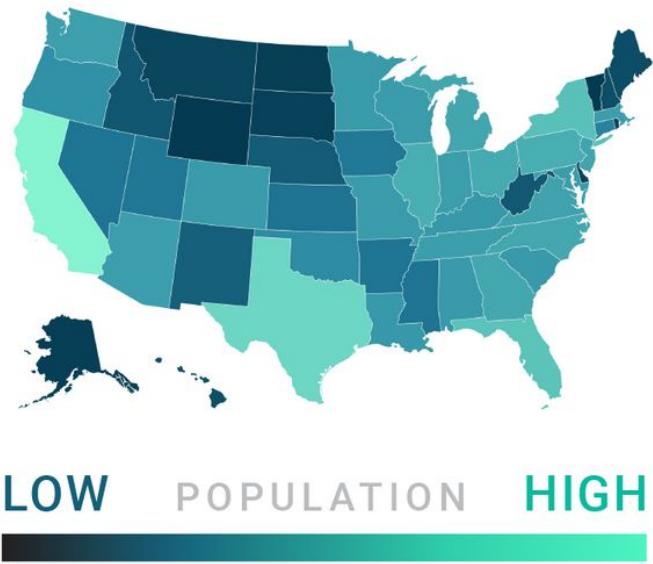
FOREST
LAKE



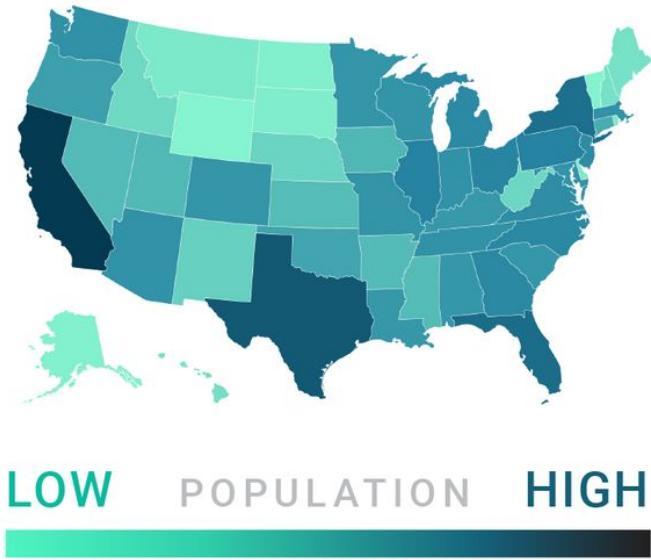
FEMALE
MALE

Use light colors for low values. Dark colors for high values.

NOT IDEAL



BETTER



Gradients are for continuous values. Distinct colors are for categories

NOT IDEAL

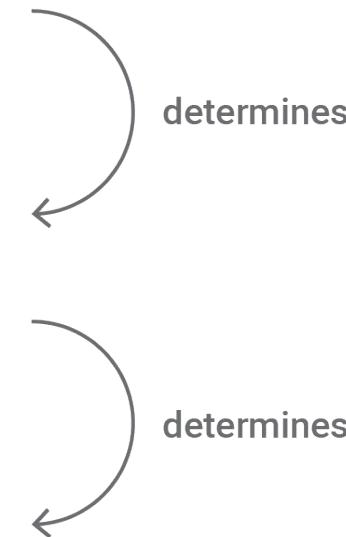
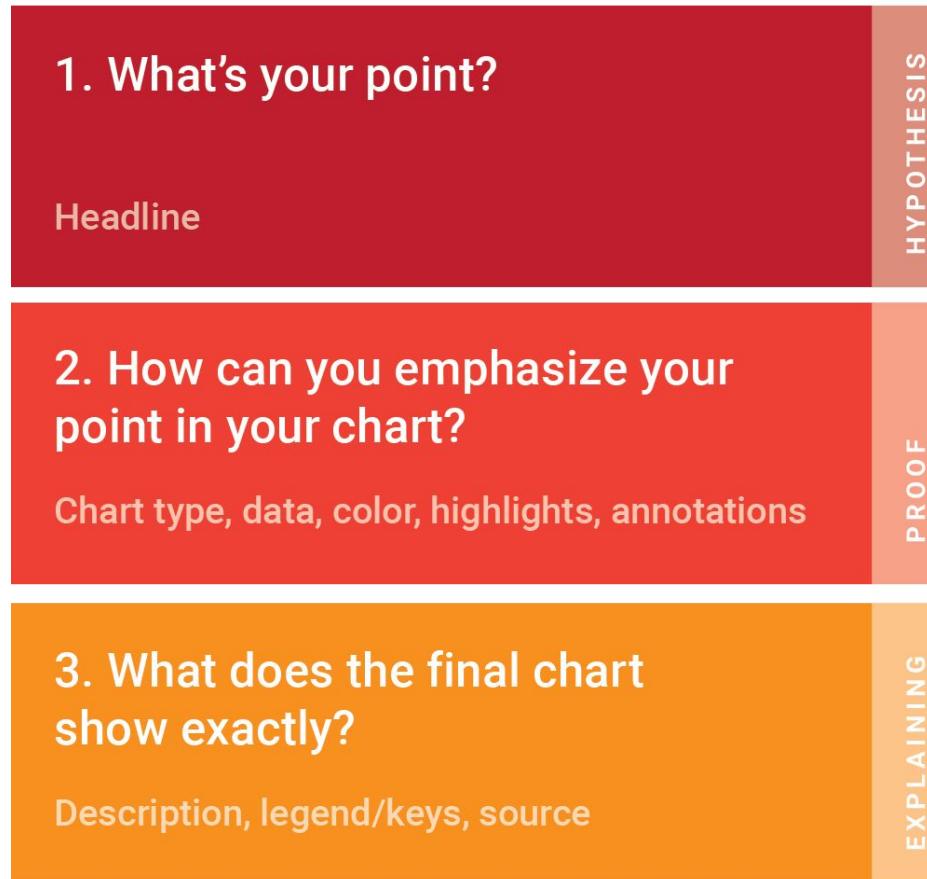


BETTER

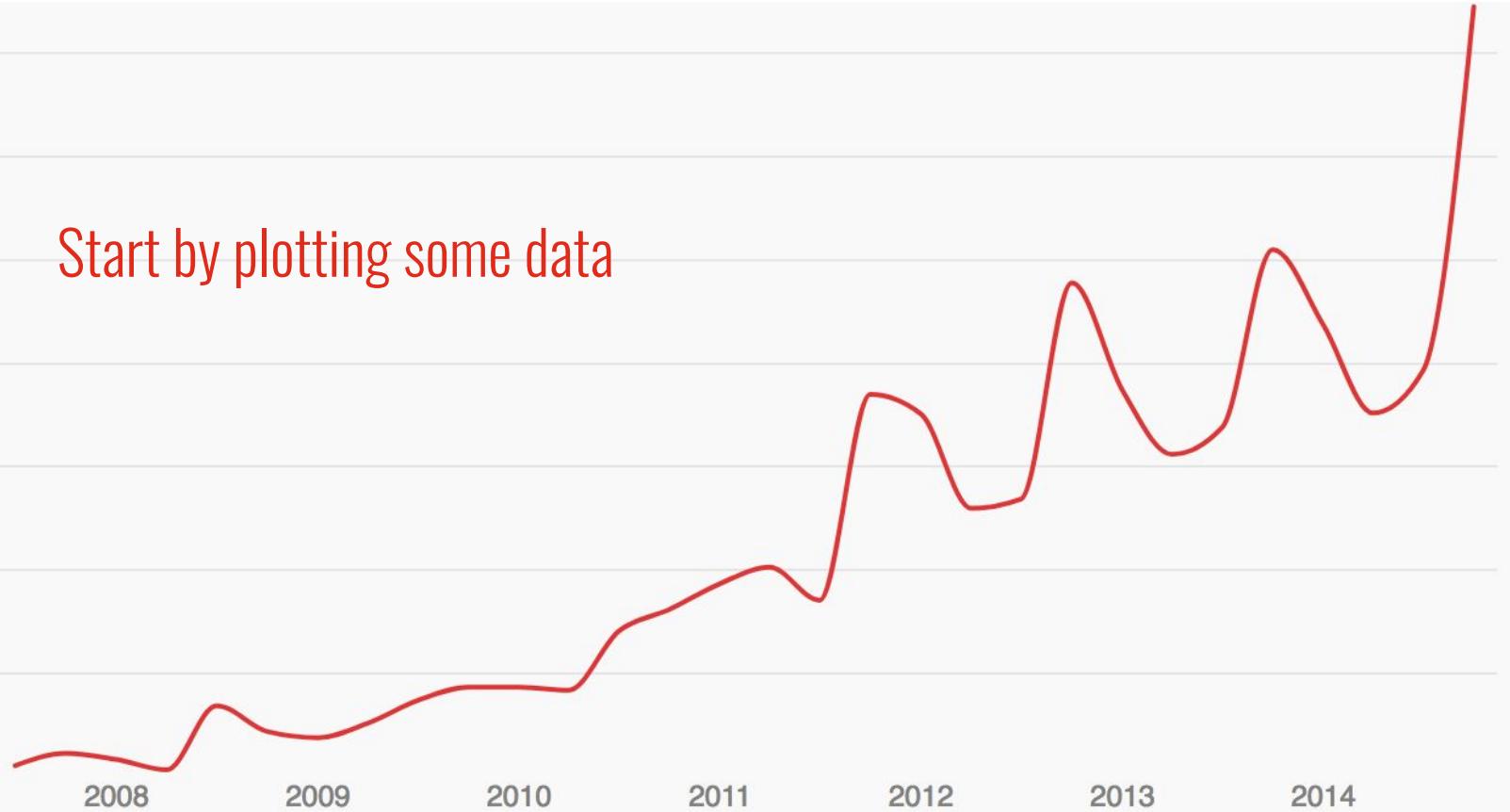


Iterative Improvement

Three questions for creating a chart

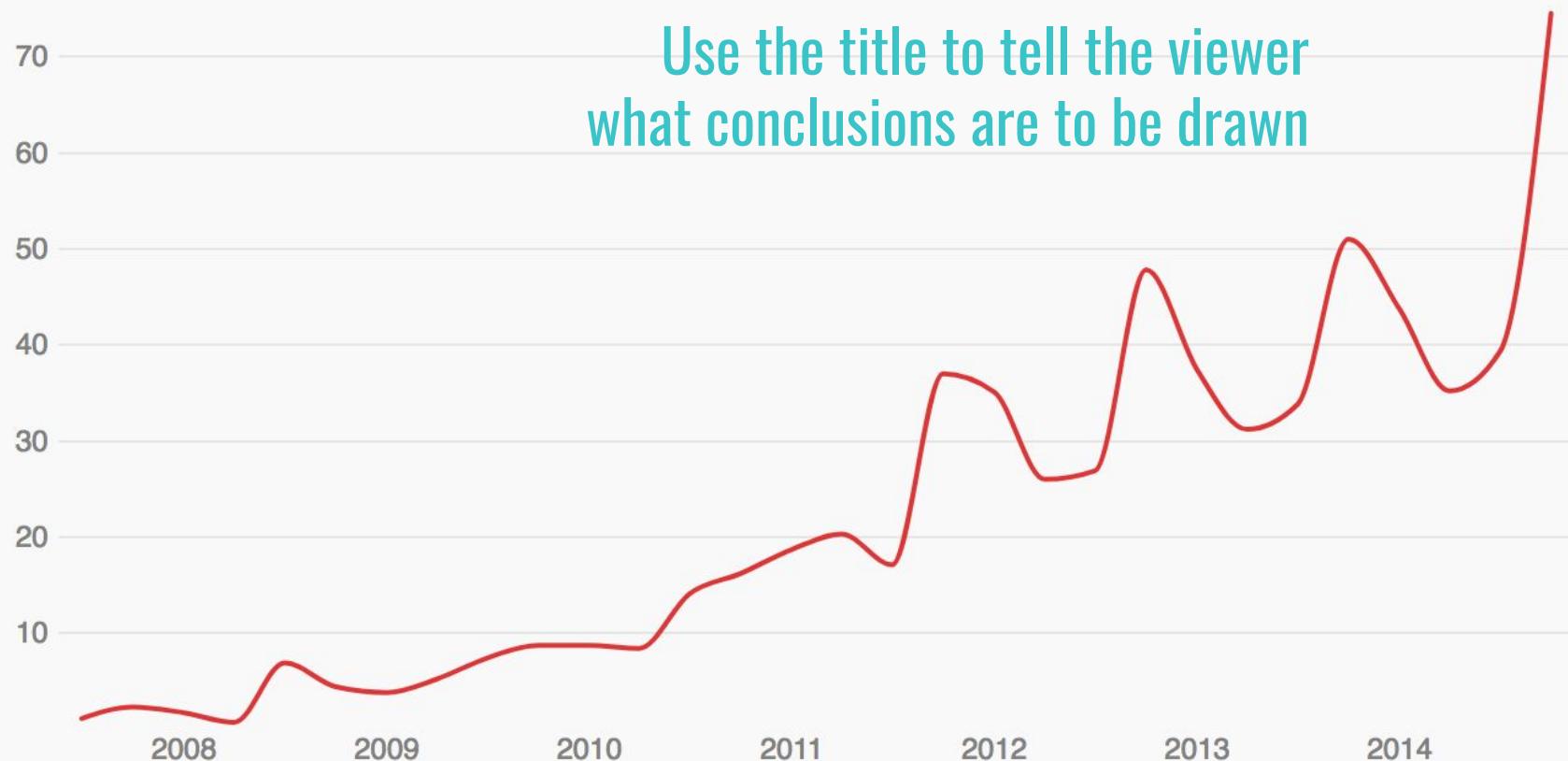


Start by plotting some data

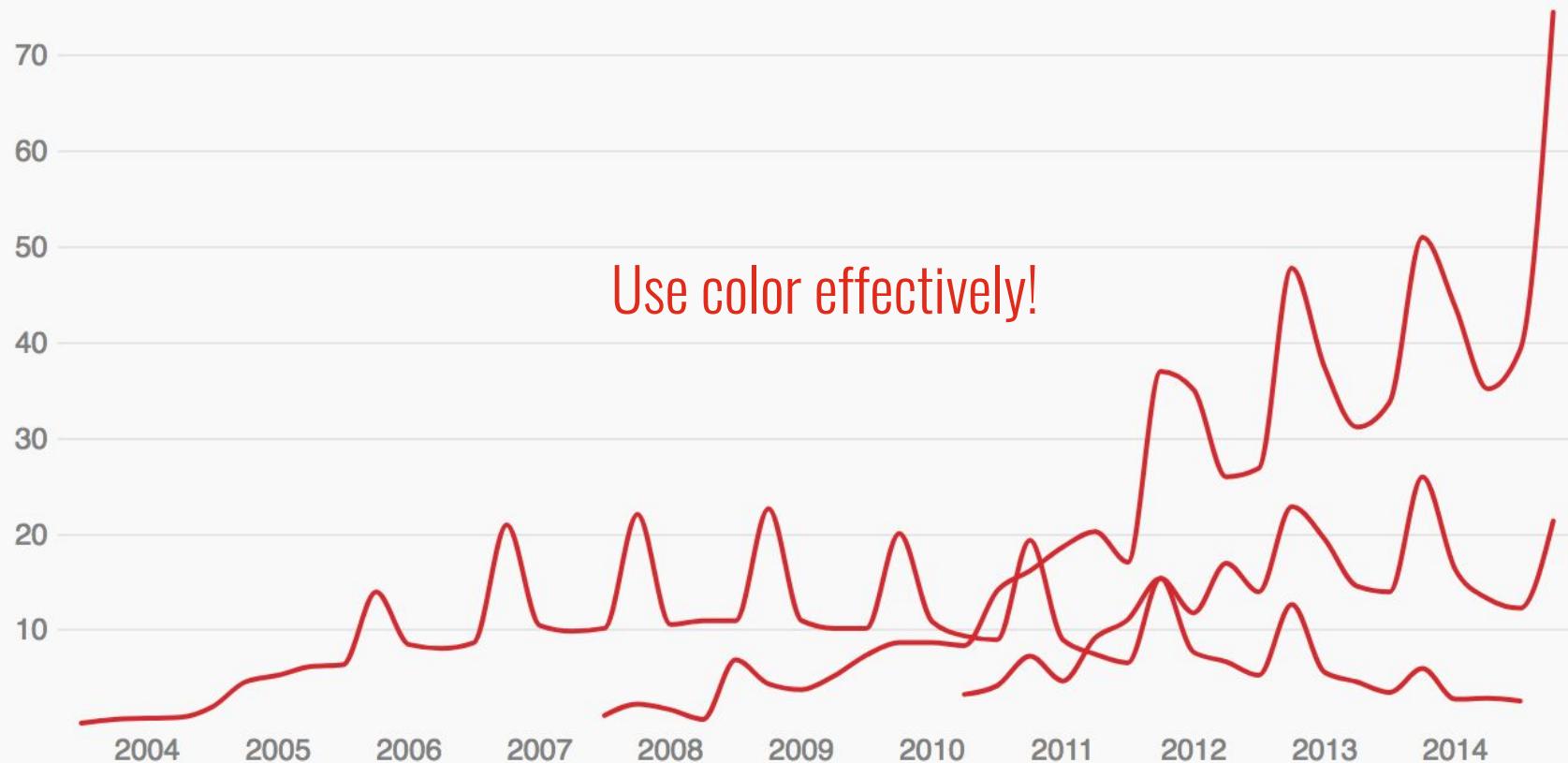


iPhone more successful than all other Apple products

Use the title to tell the viewer
what conclusions are to be drawn

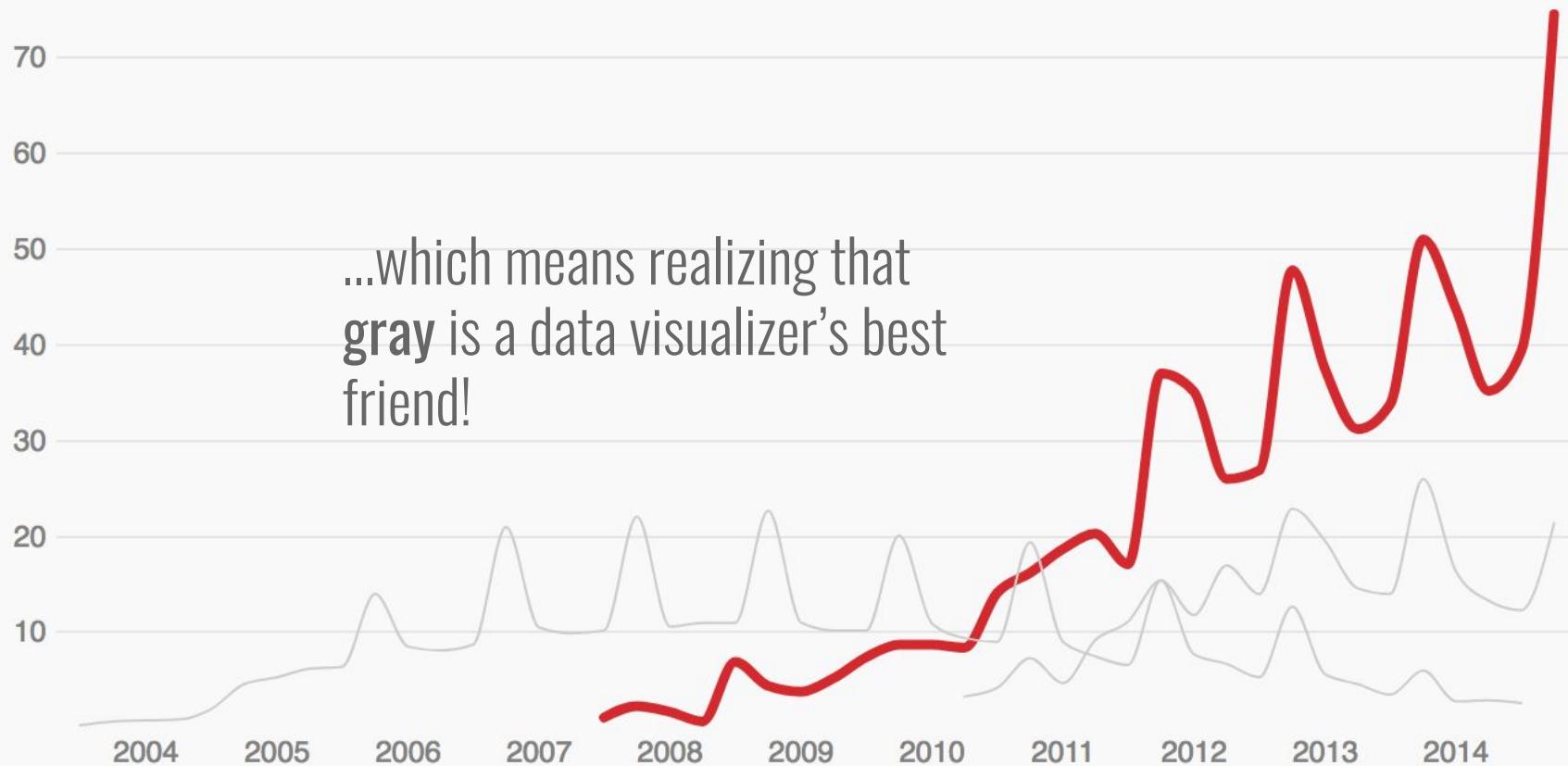


iPhone more successful than all other Apple products



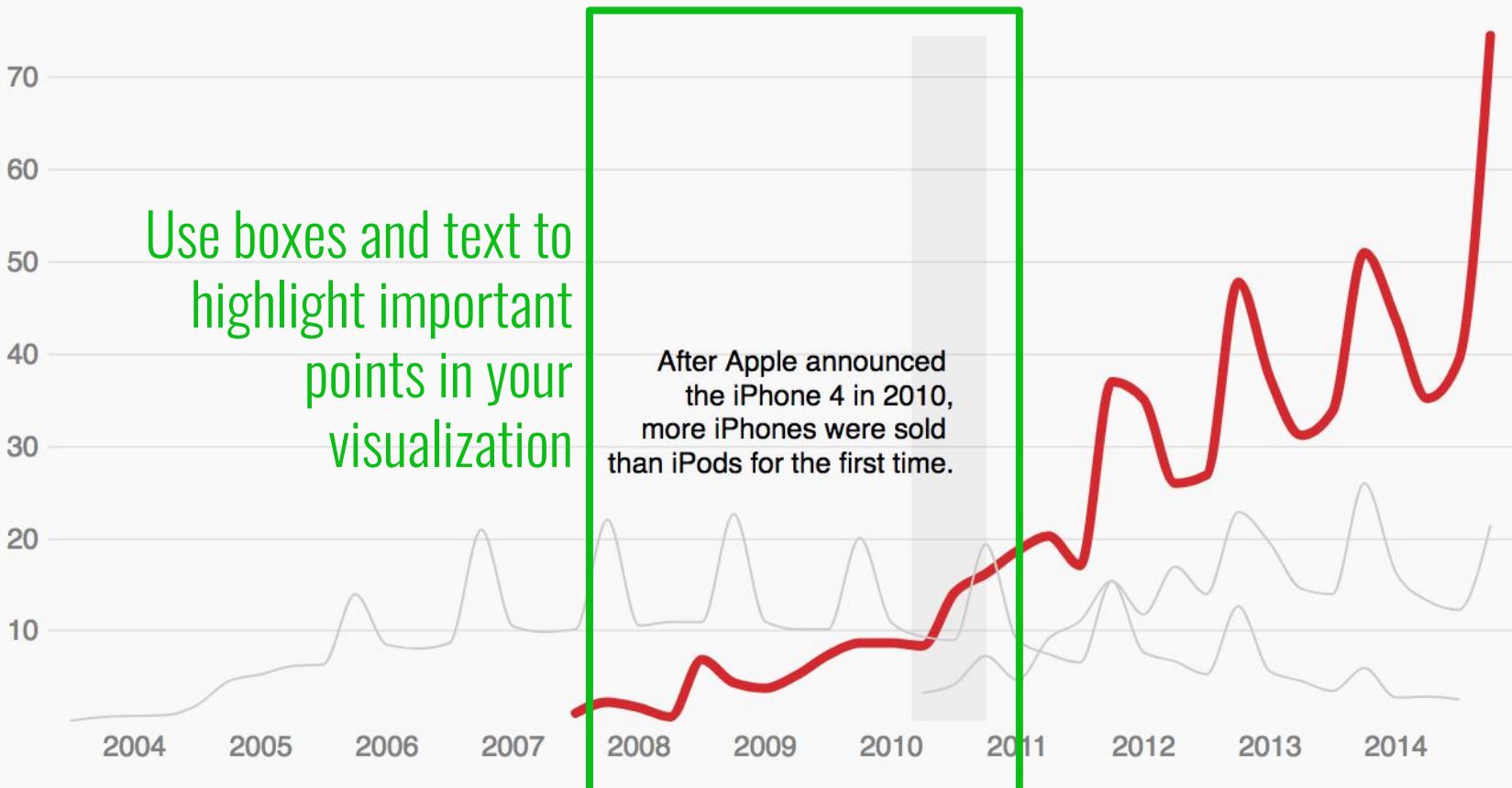
iPhone more successful than all other Apple products

...which means realizing that
gray is a data visualizer's best
friend!



iPhone more successful than all other Apple products

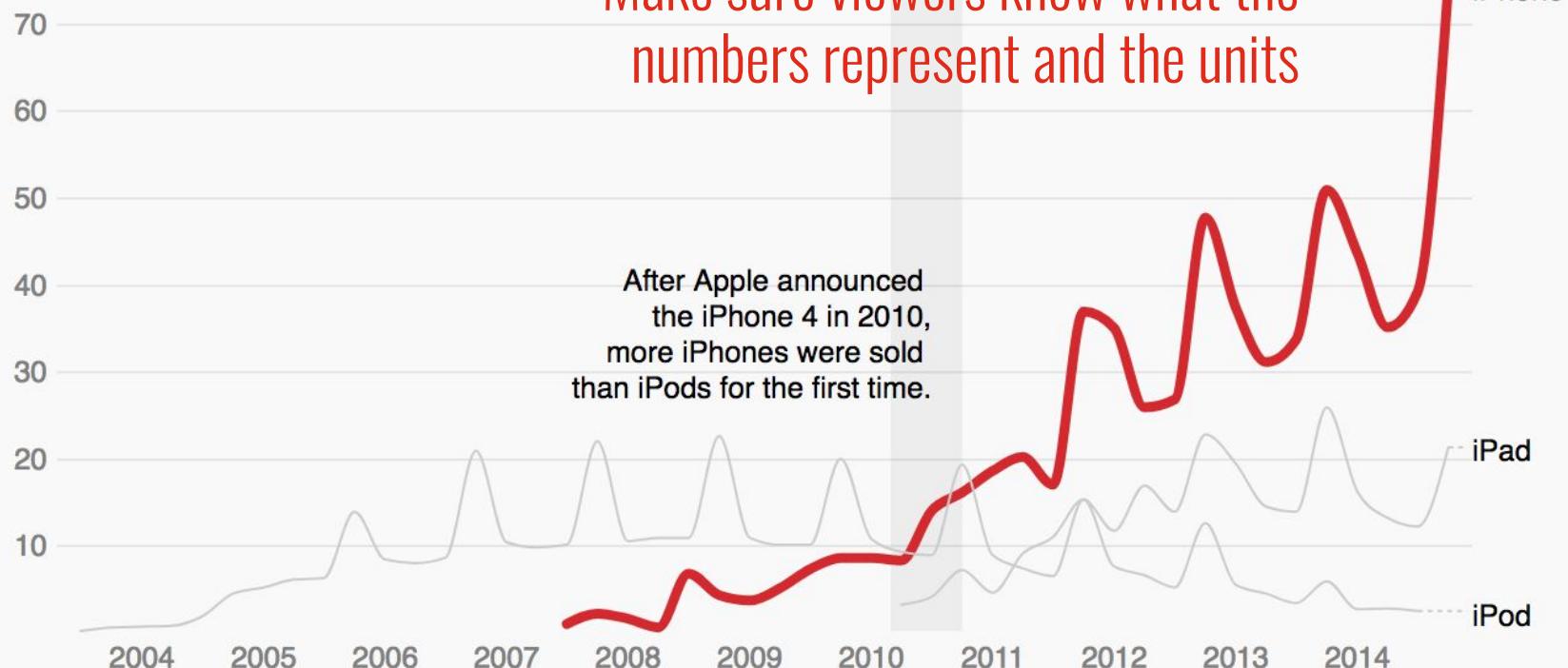
Use boxes and text to highlight important points in your visualization



iPhone more successful than all other Apple products

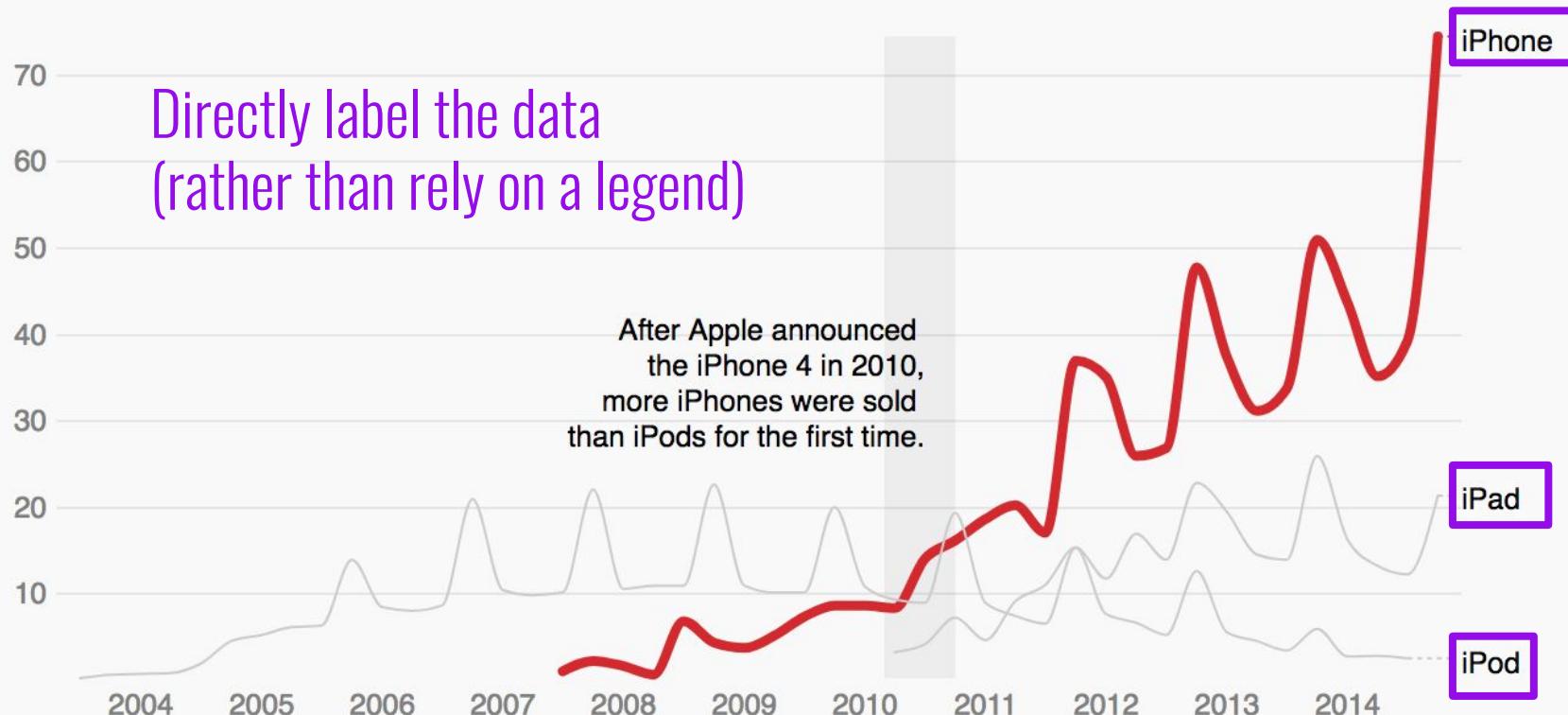
Worldwide sales of selected Apple products in million, by fiscal quarter, 2000 to 2014

Make sure viewers know what the numbers represent and the units



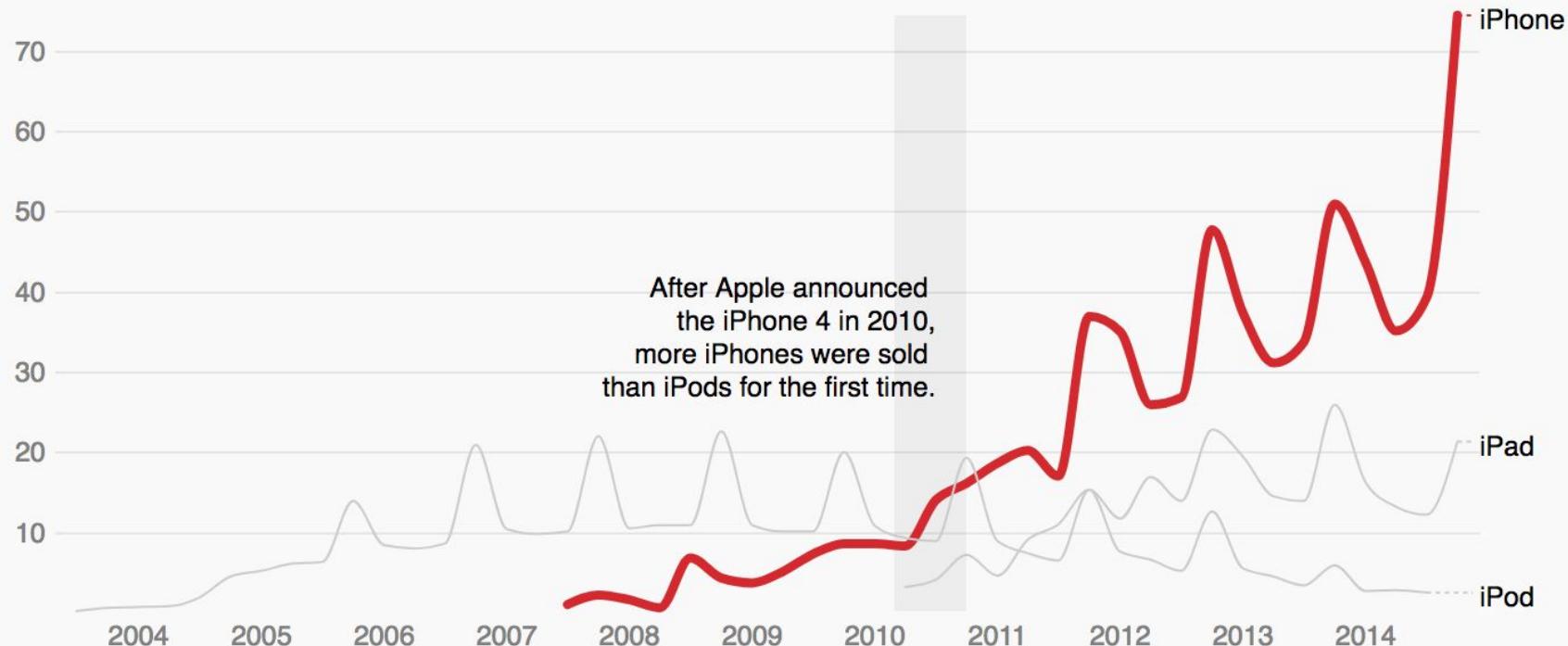
iPhone more successful than all other Apple products

Worldwide sales of selected Apple products in million, by fiscal quarter, 2000 to 2014



iPhone more successful than all other Apple products

Worldwide sales of selected Apple products in million, by fiscal quarter, 2000 to 2014

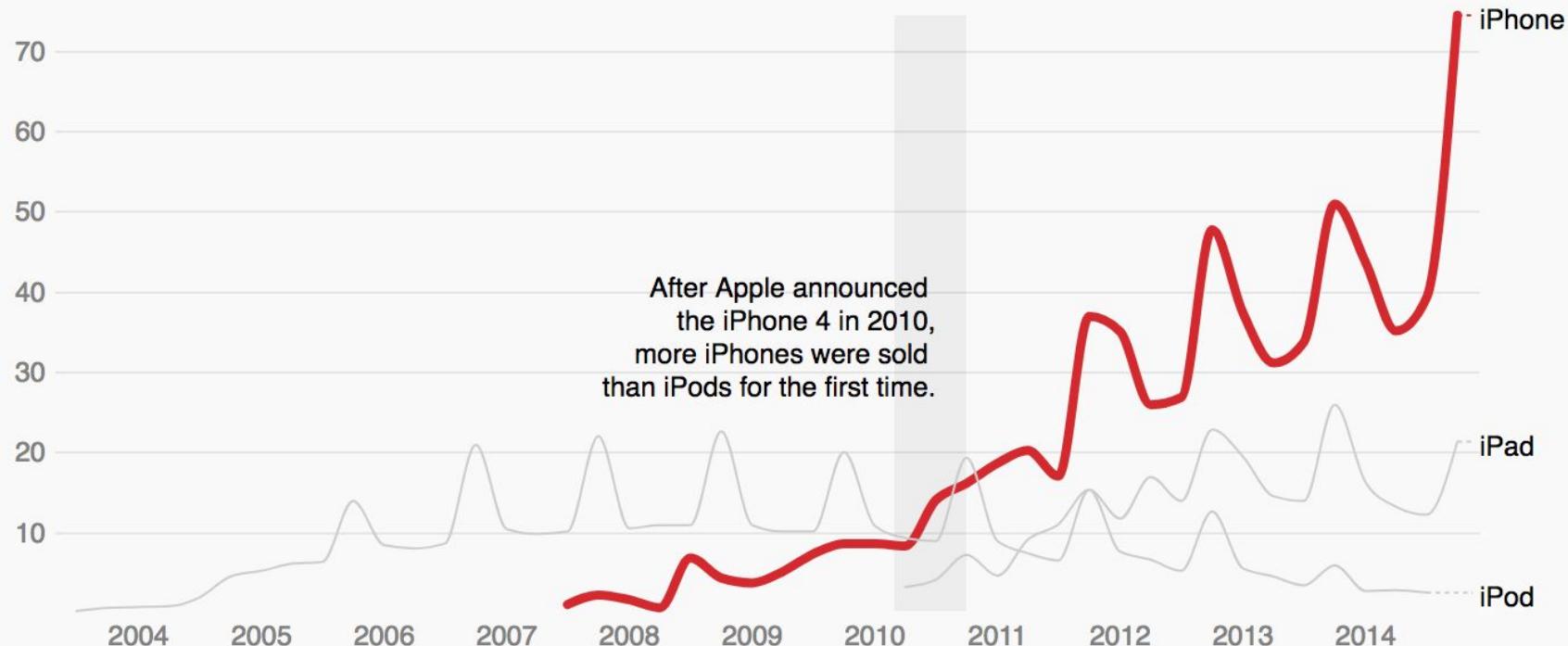


Source: Apple Inc. • Get the data • Created with Datawrapper

Always include your source

iPhone more successful than all other Apple products

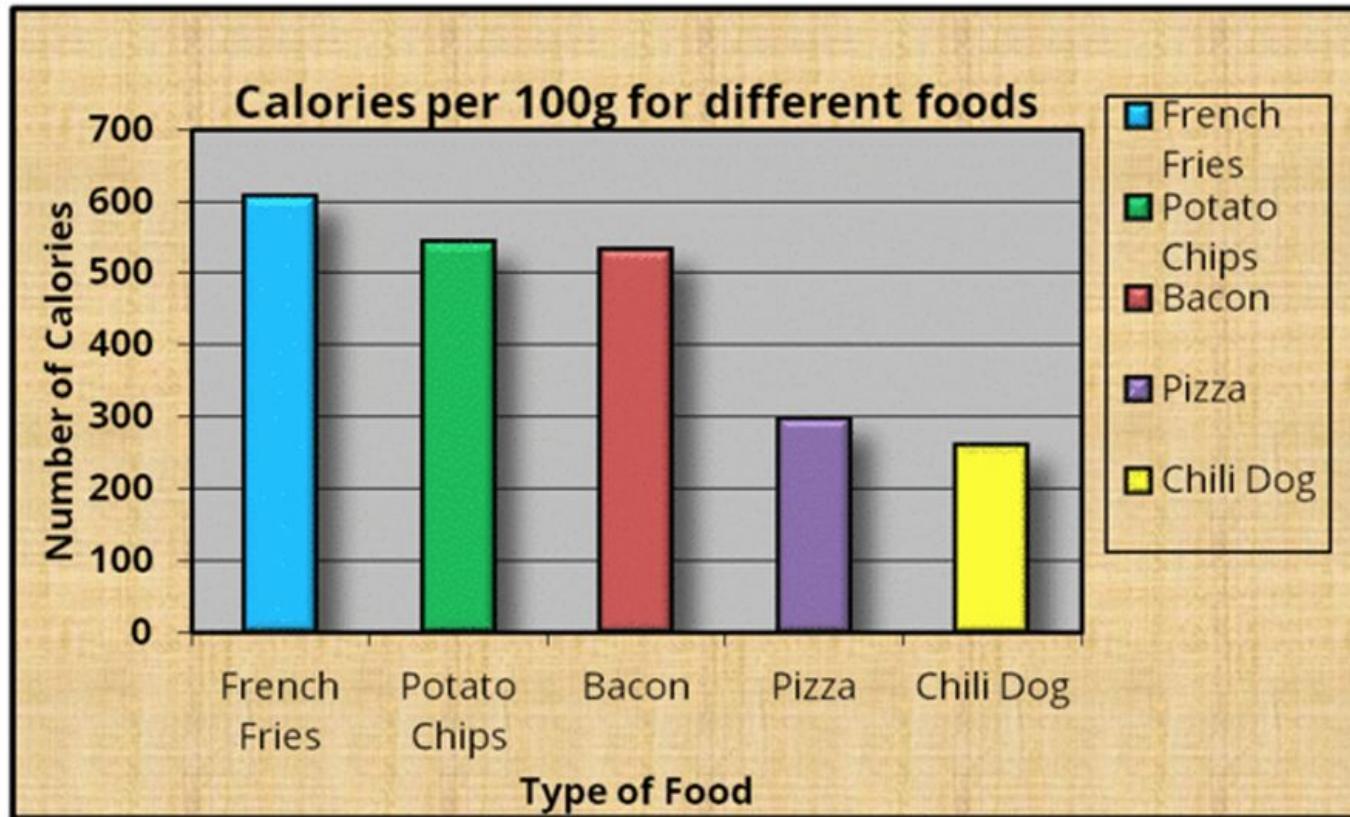
Worldwide sales of selected Apple products in million, by fiscal quarter, 2000 to 2014



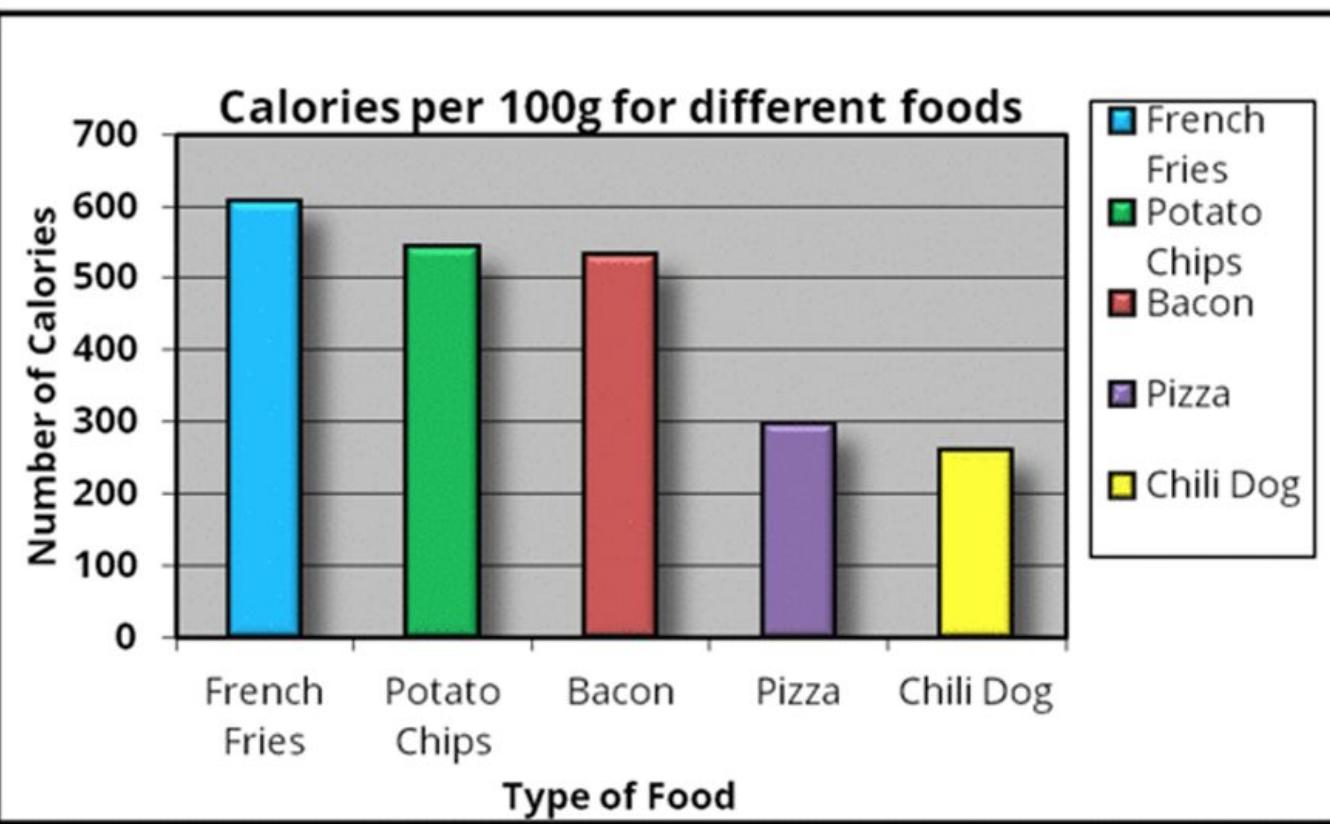
Source: Apple Inc. • Get the data • Created with Datawrapper

AIM TO IMPROVE YOUR:
**data:ink
ratio**

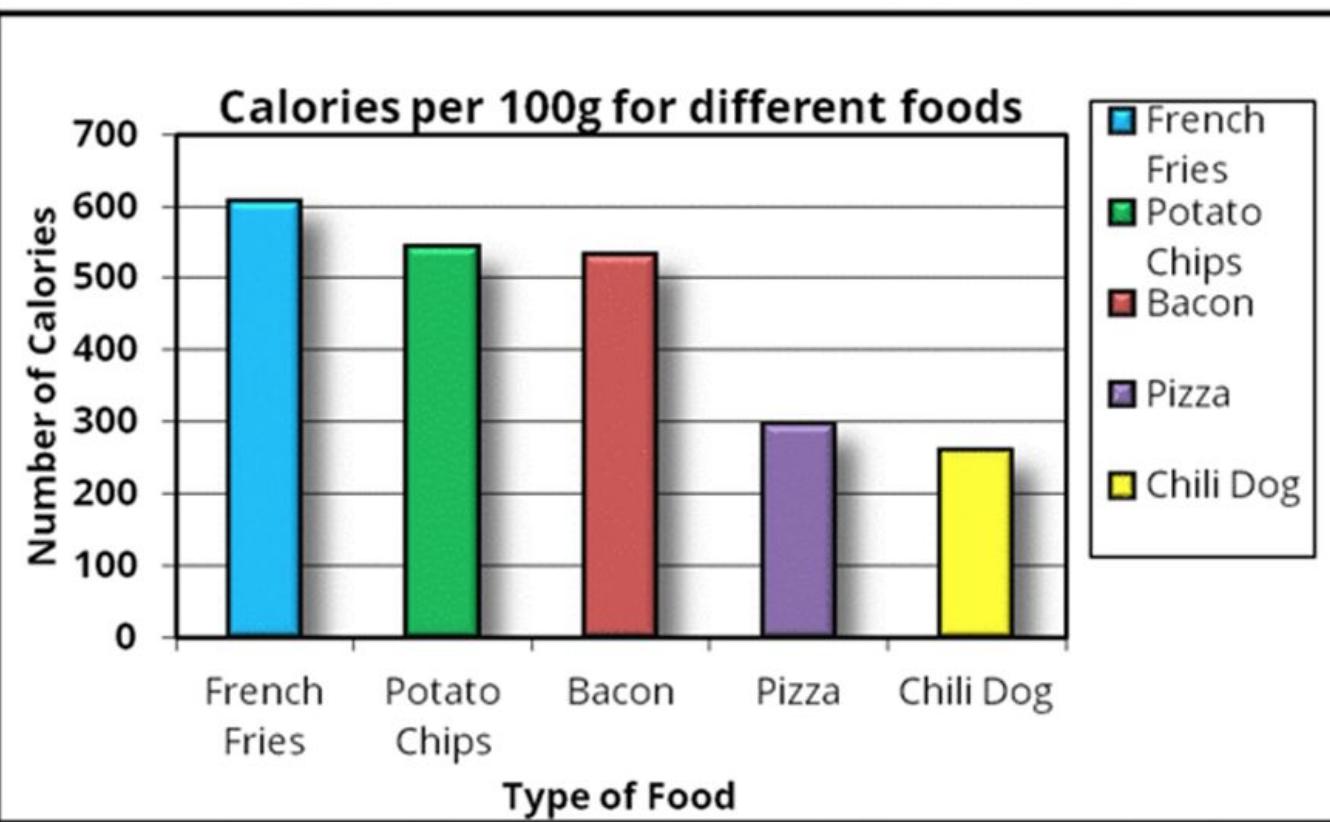
Remove backgrounds



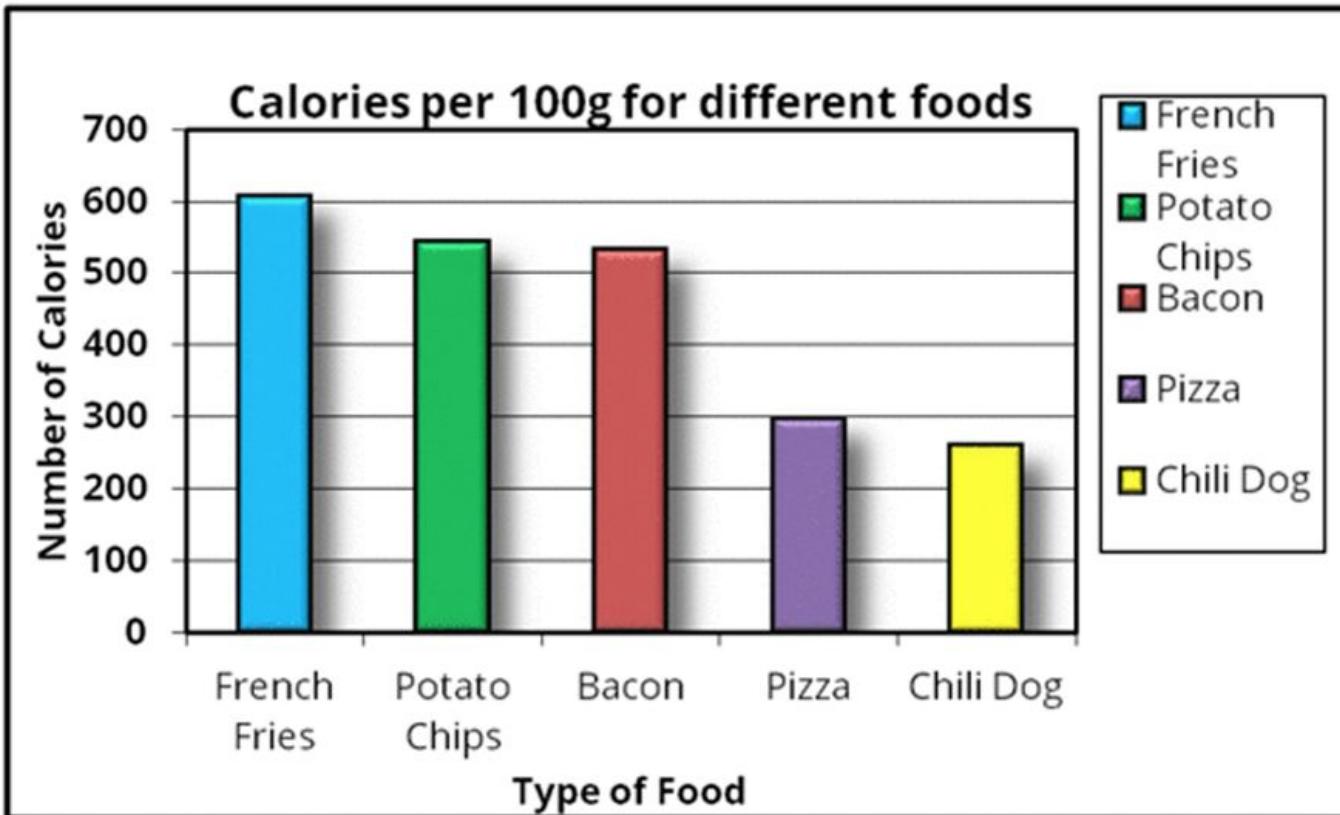
Remove backgrounds



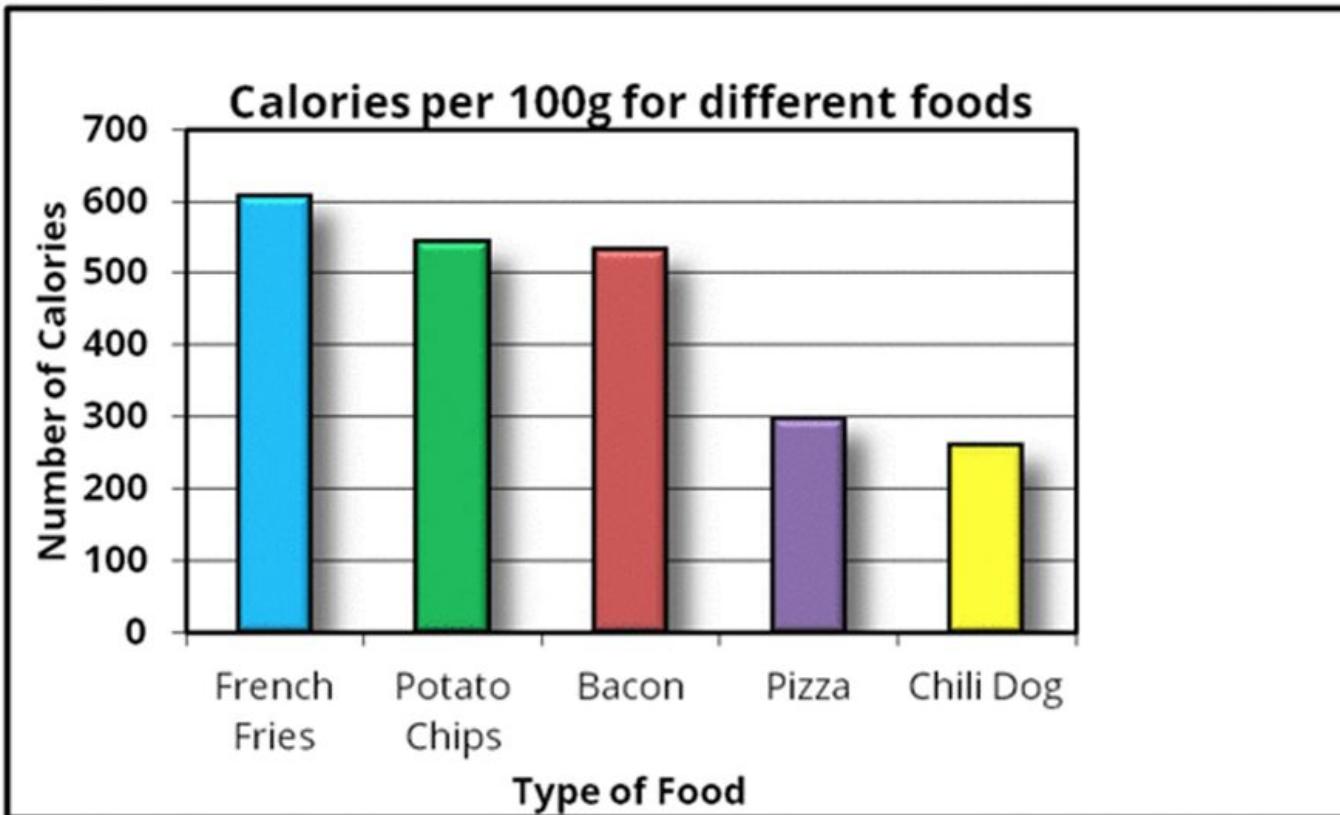
Remove backgrounds



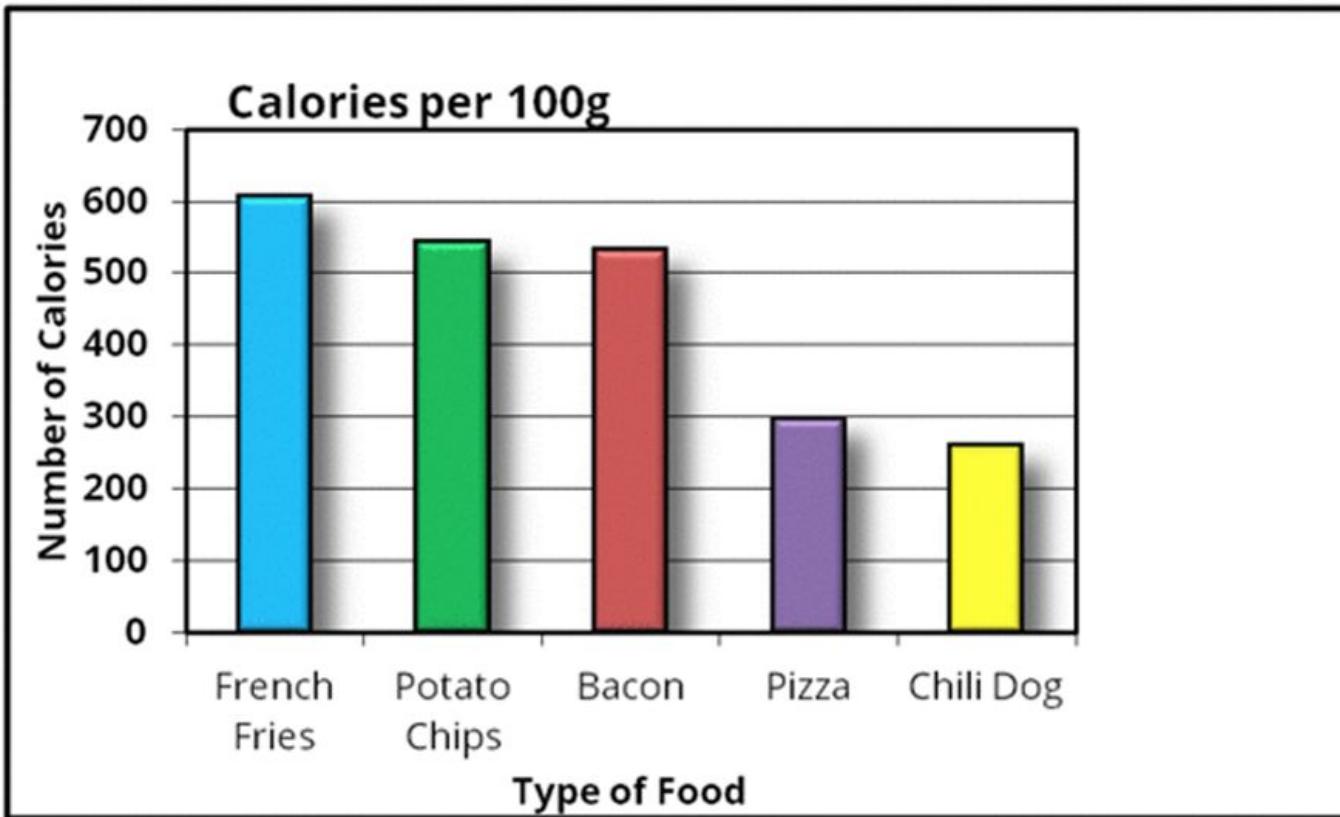
Remove redundant labels



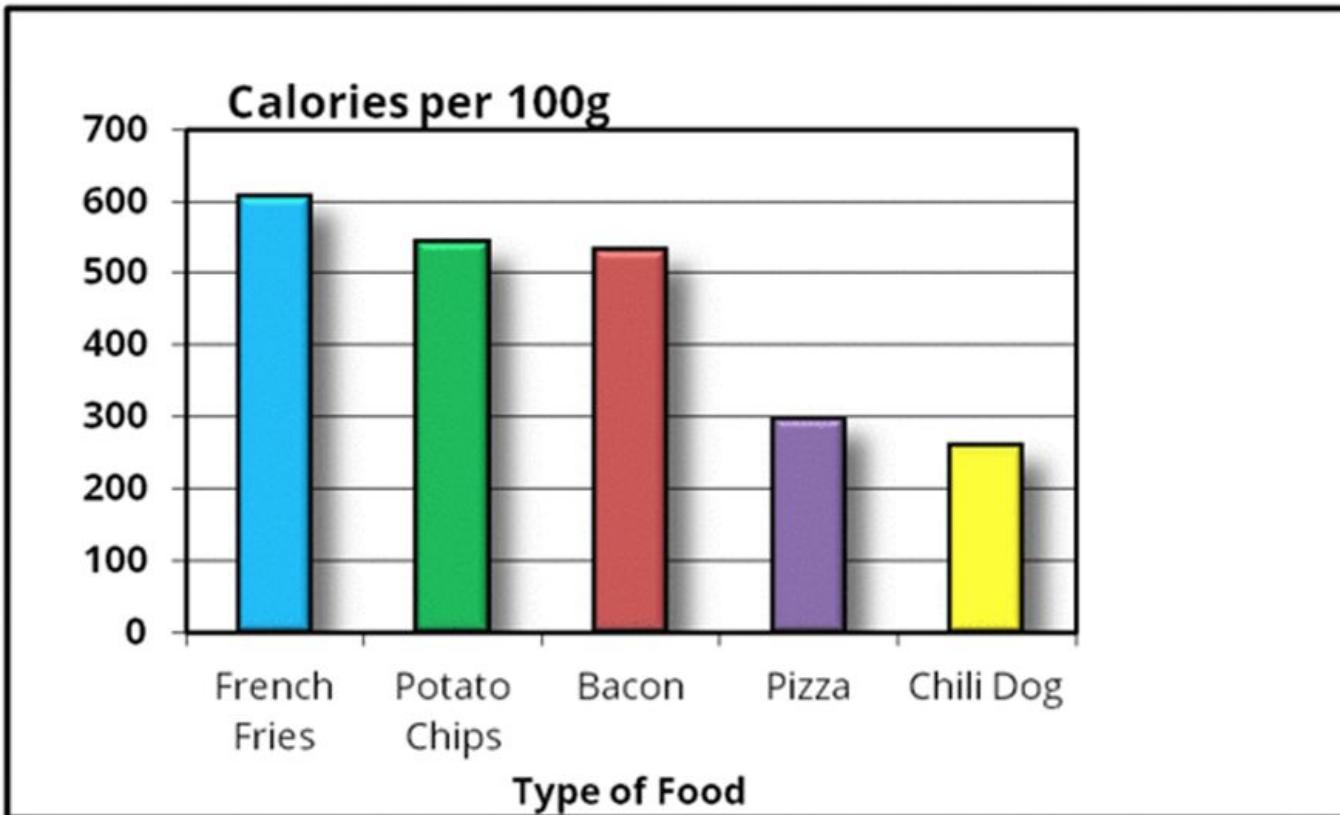
Remove redundant labels



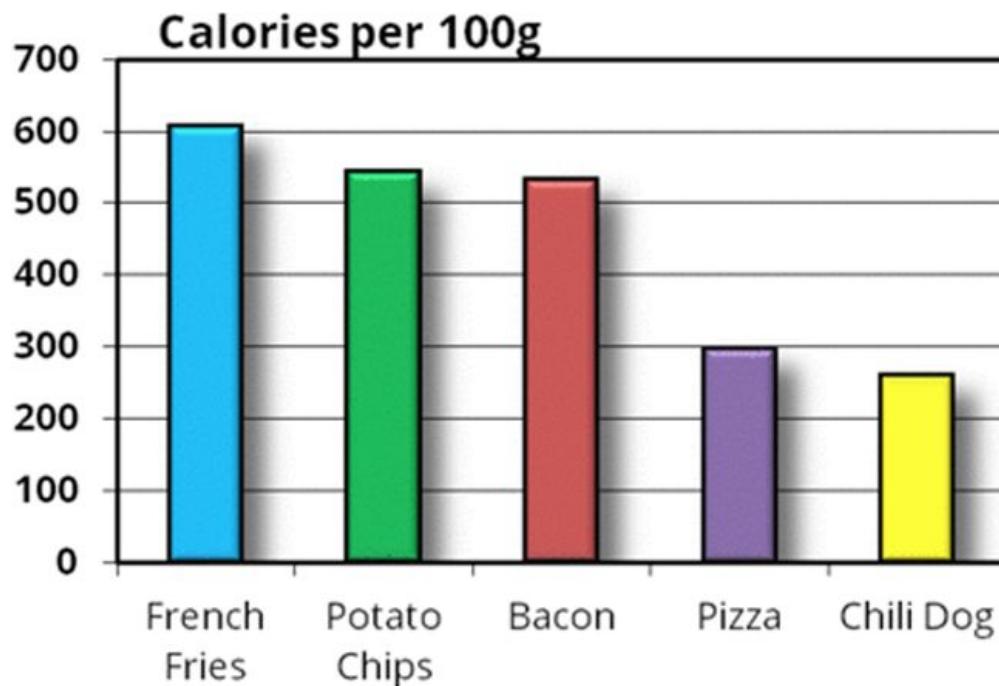
Remove redundant labels



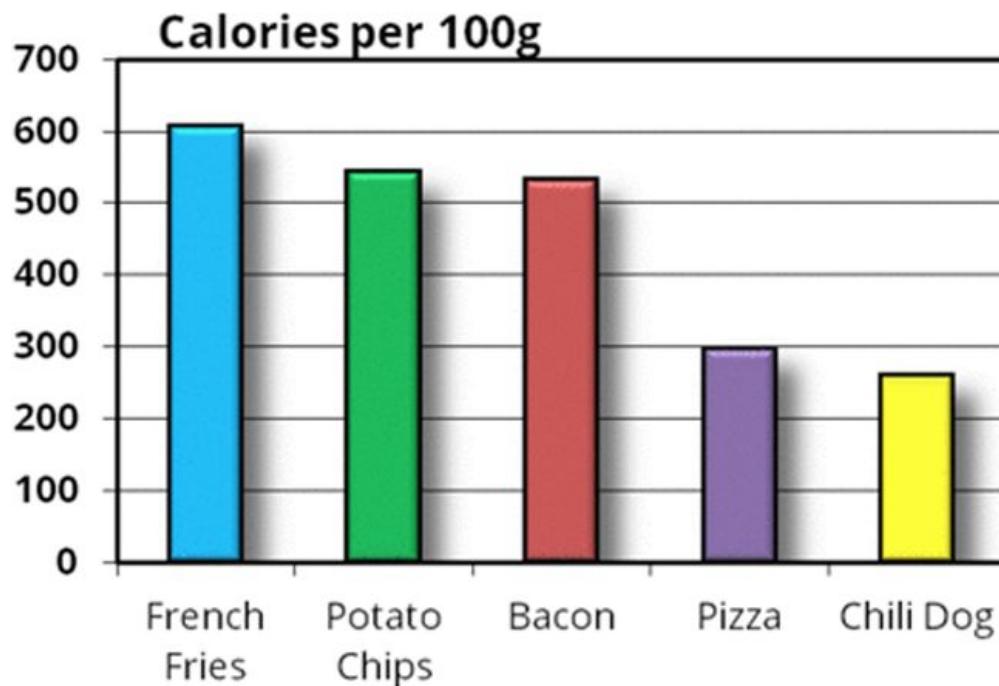
Remove redundant labels



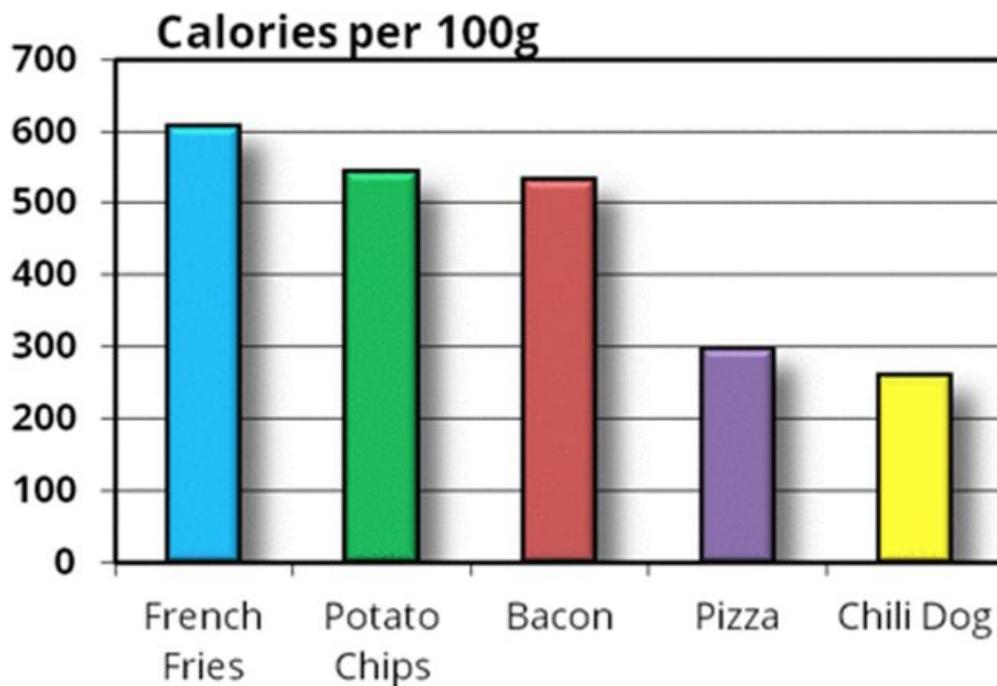
Remove redundant labels



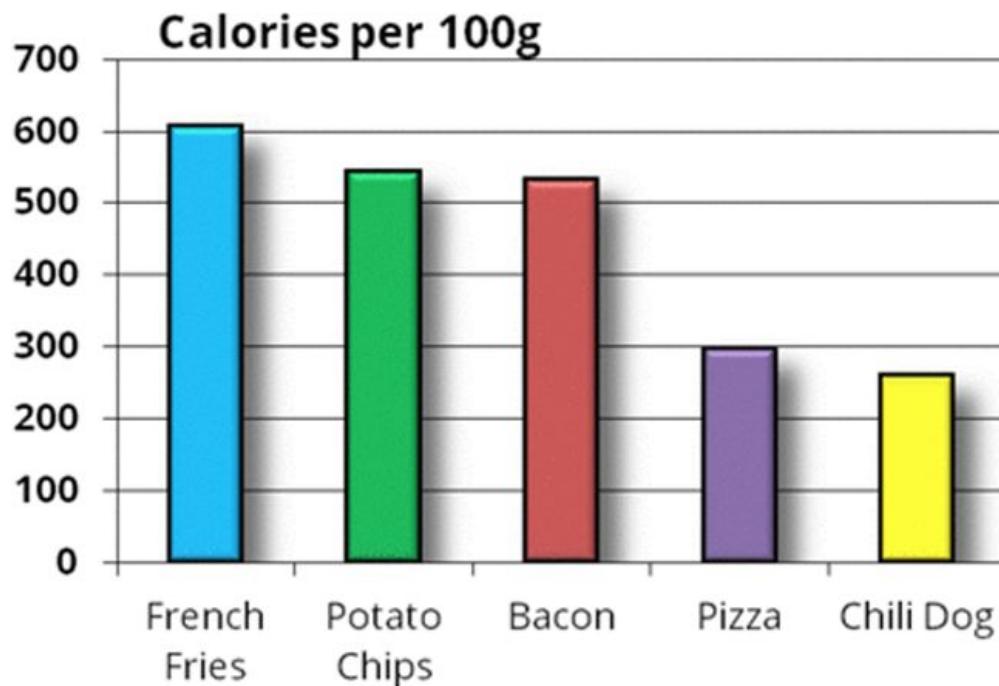
Remove borders



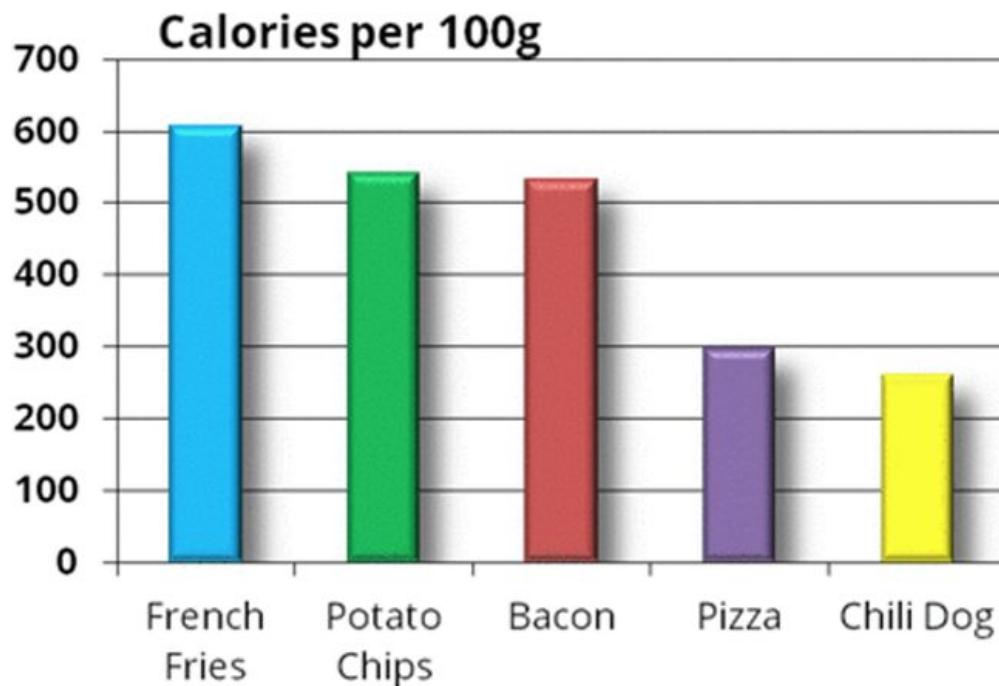
Remove borders



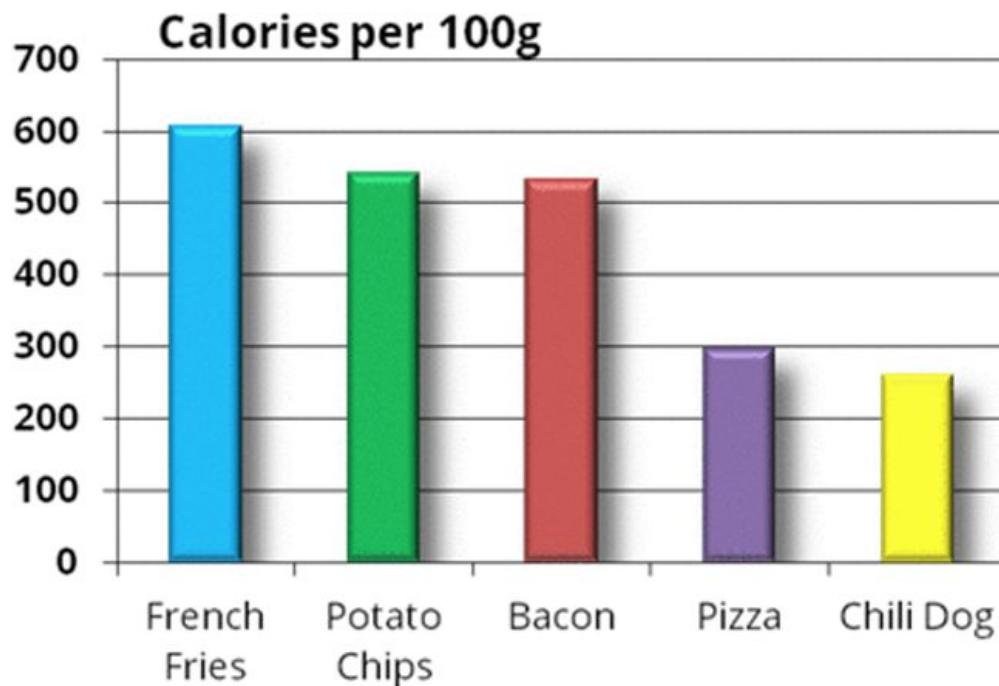
Remove borders



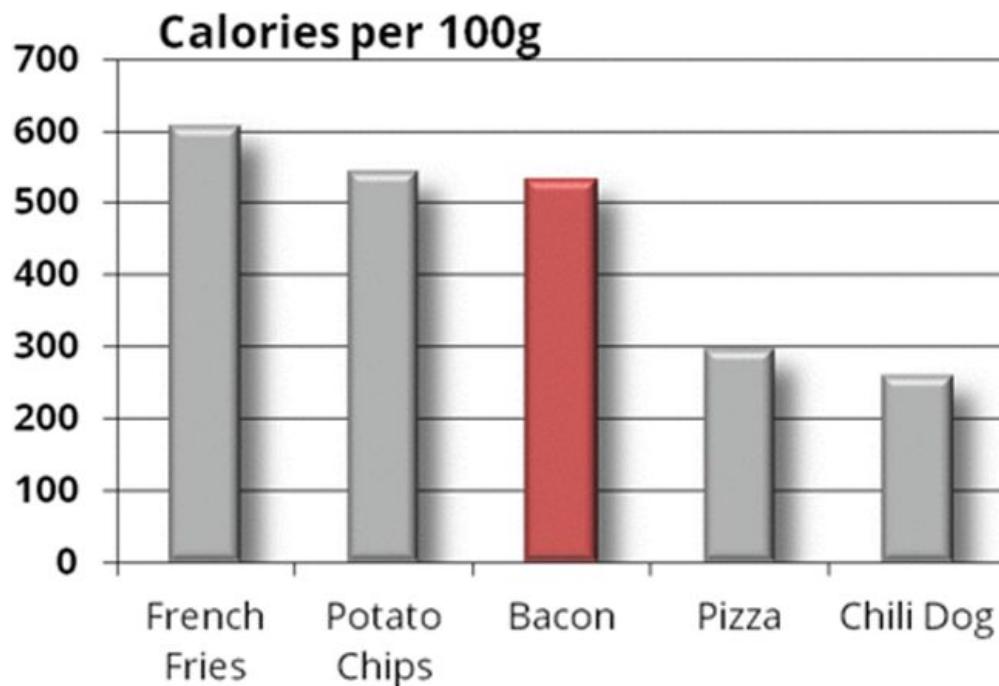
Remove borders



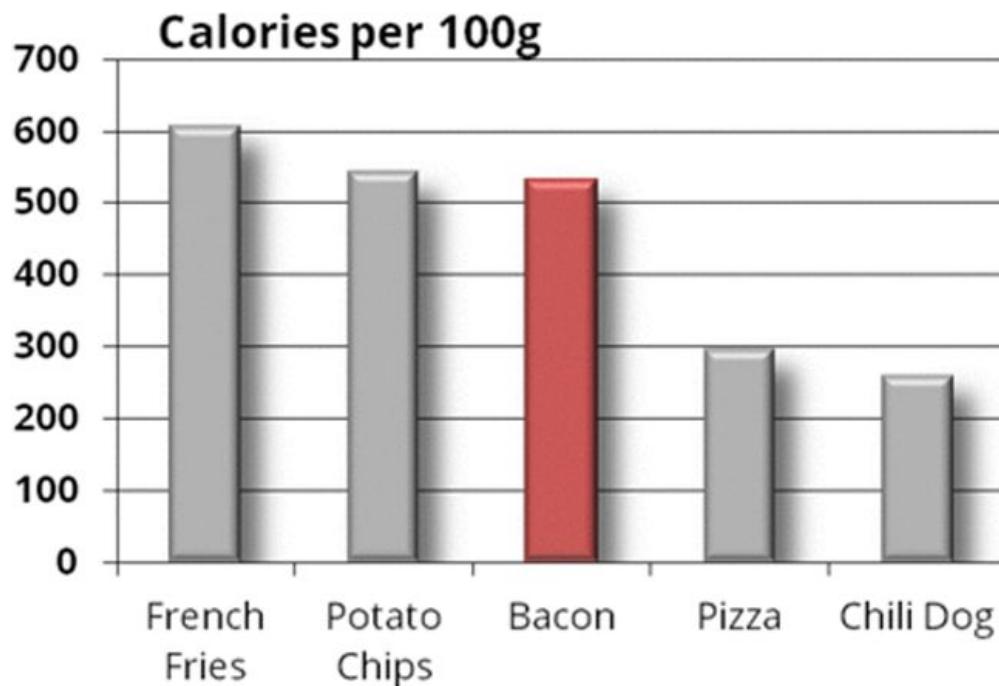
Reduce colors



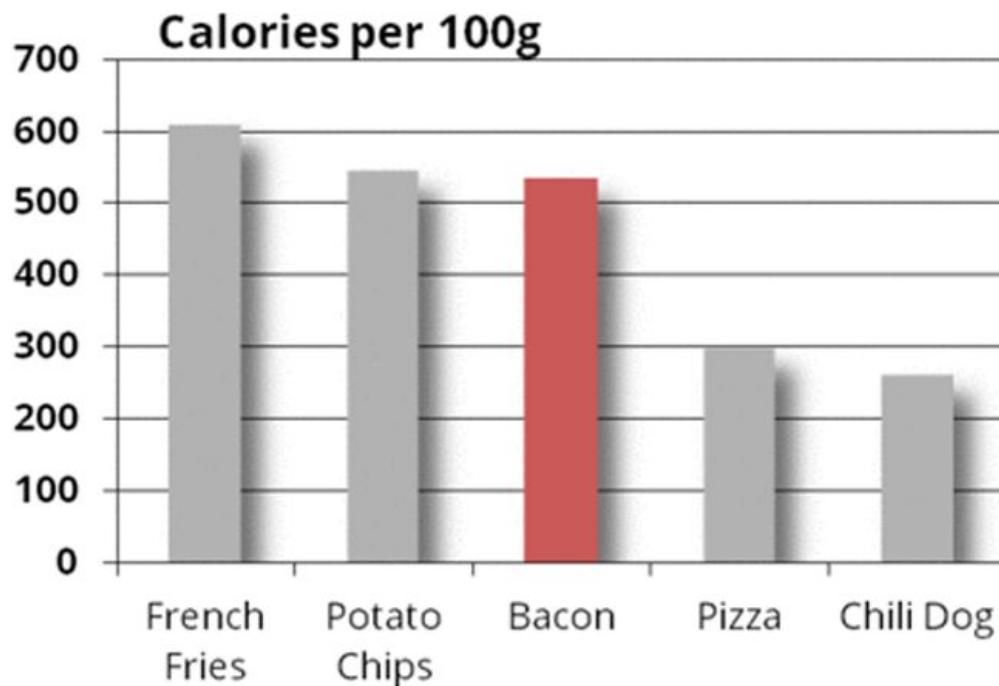
Reduce colors



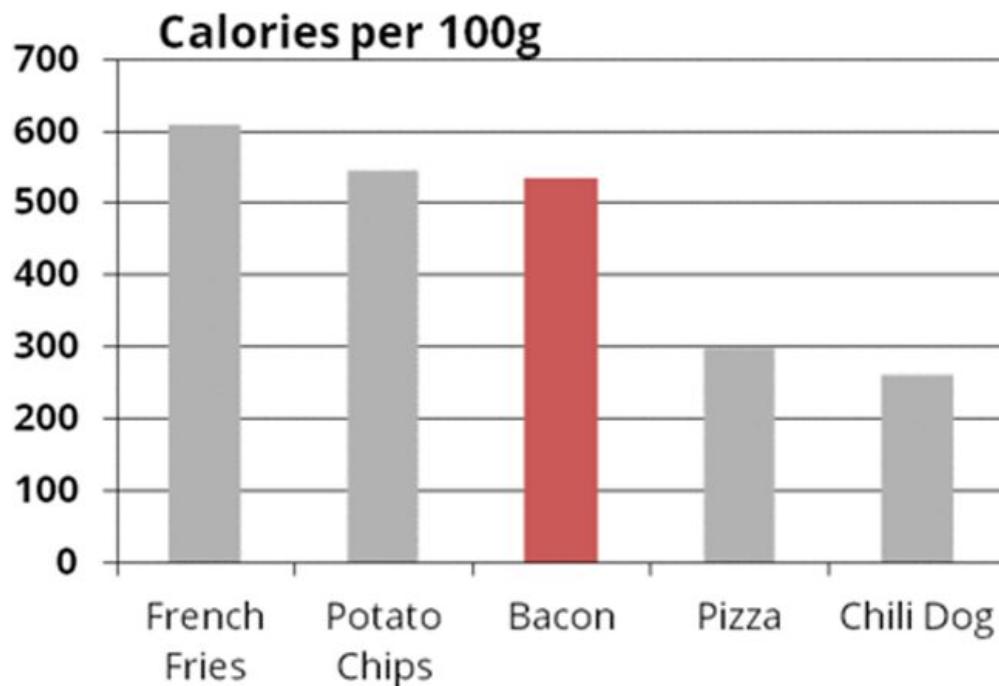
Remove special effects



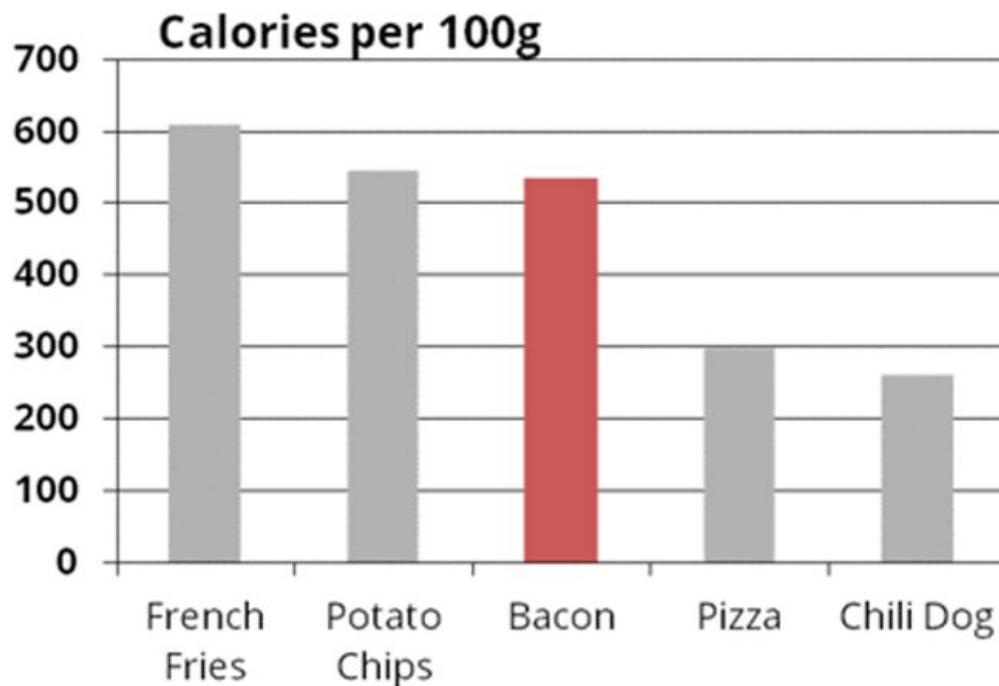
Remove special effects



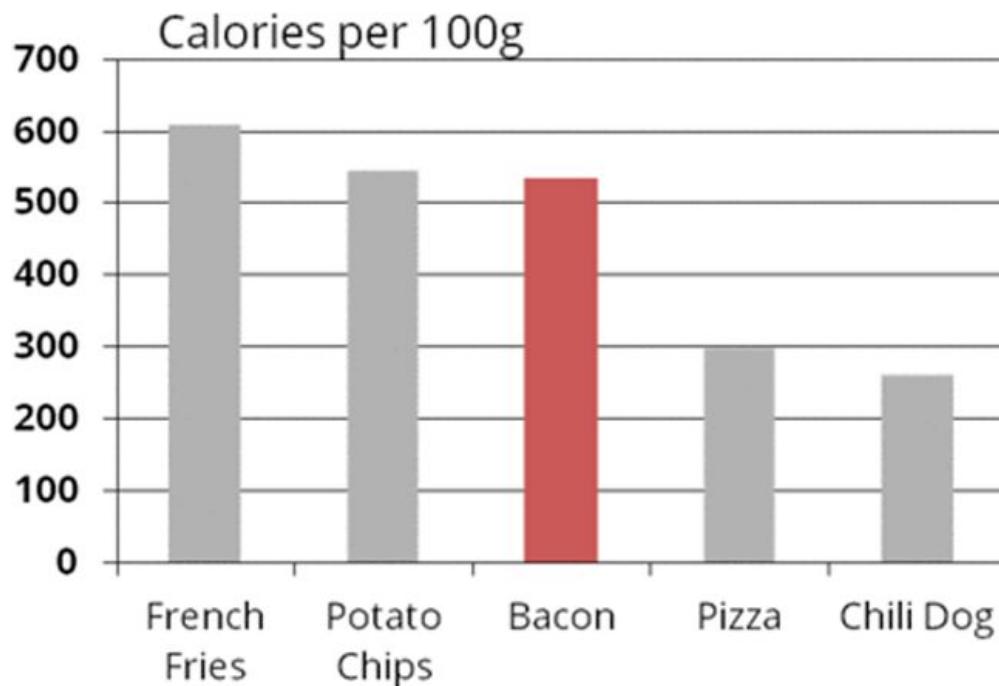
Remove special effects



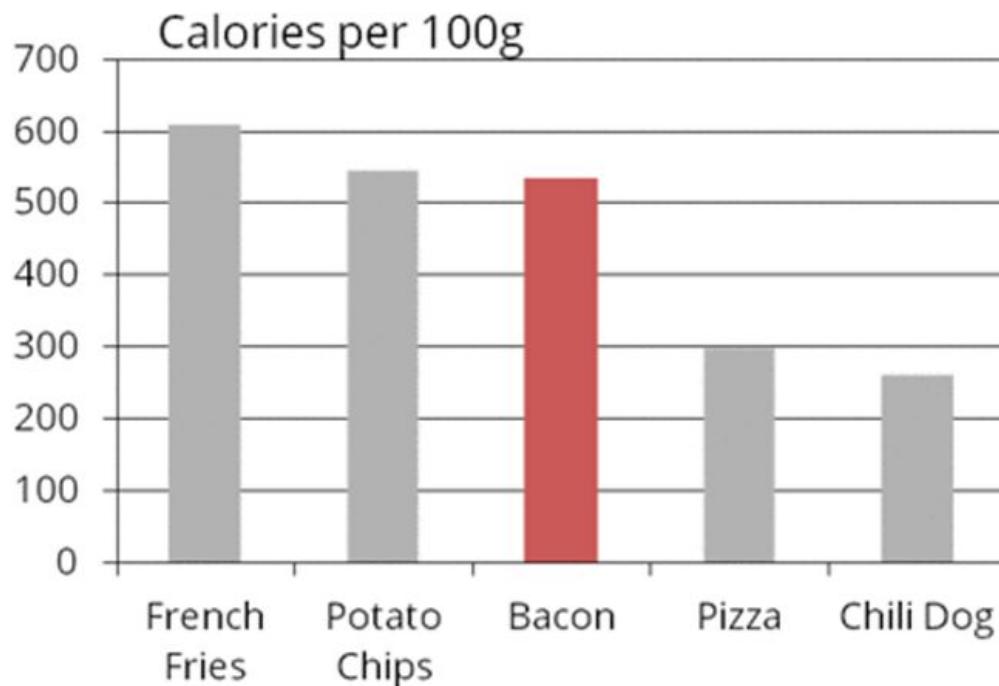
Remove bolding



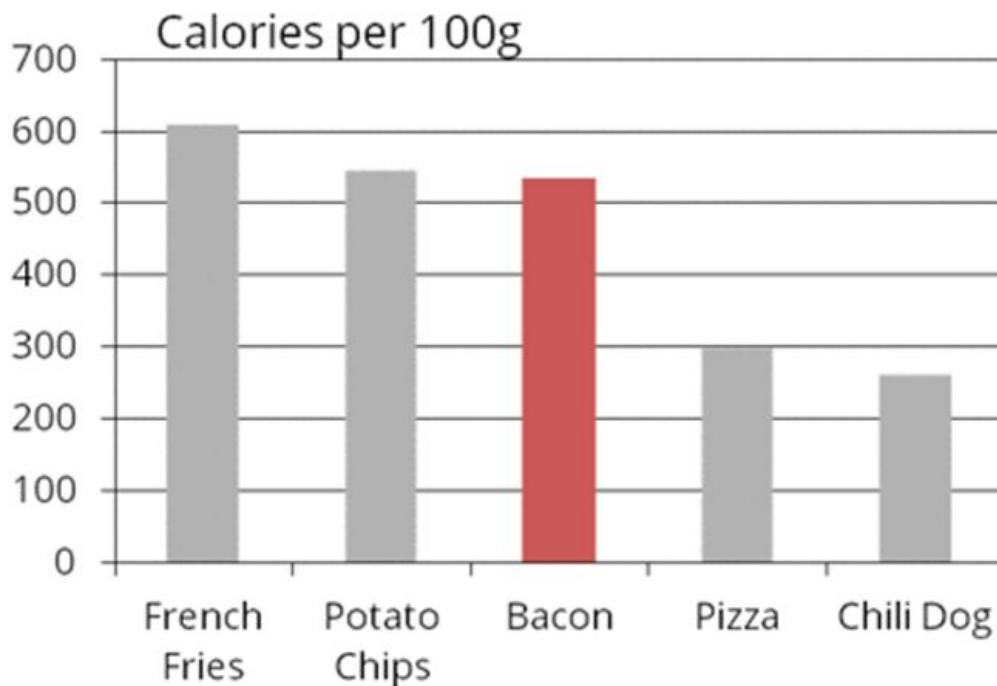
Remove bolding



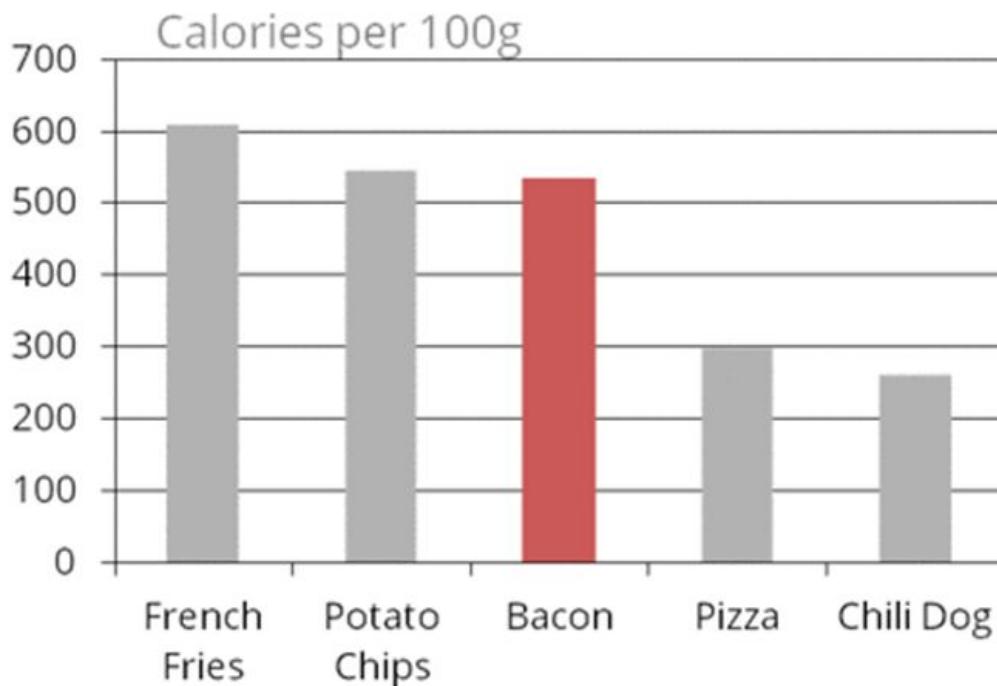
Remove bolding



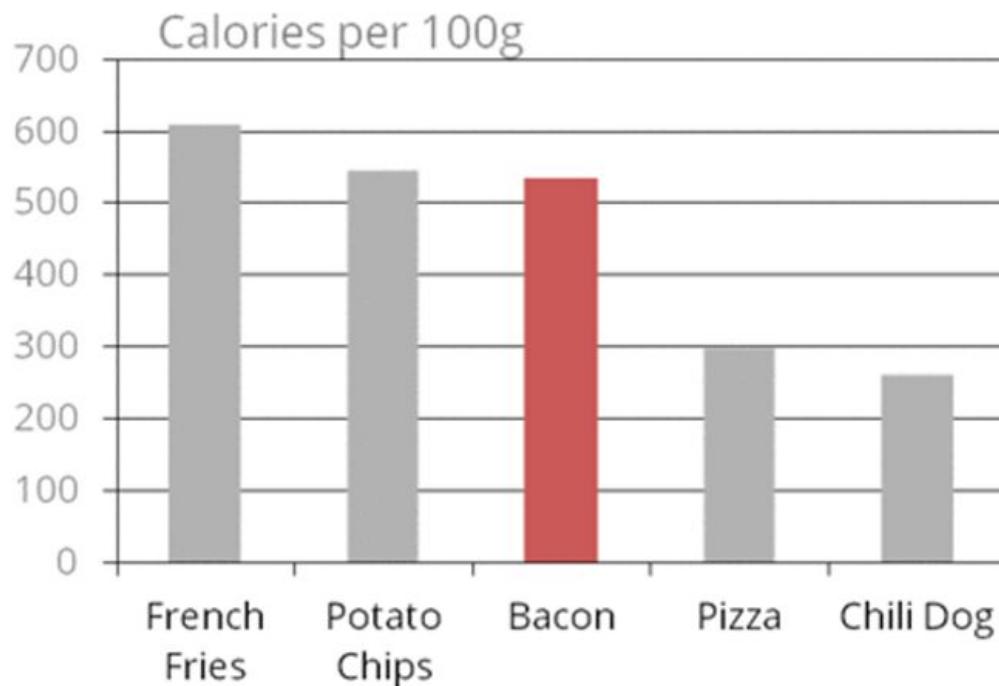
Lighten labels



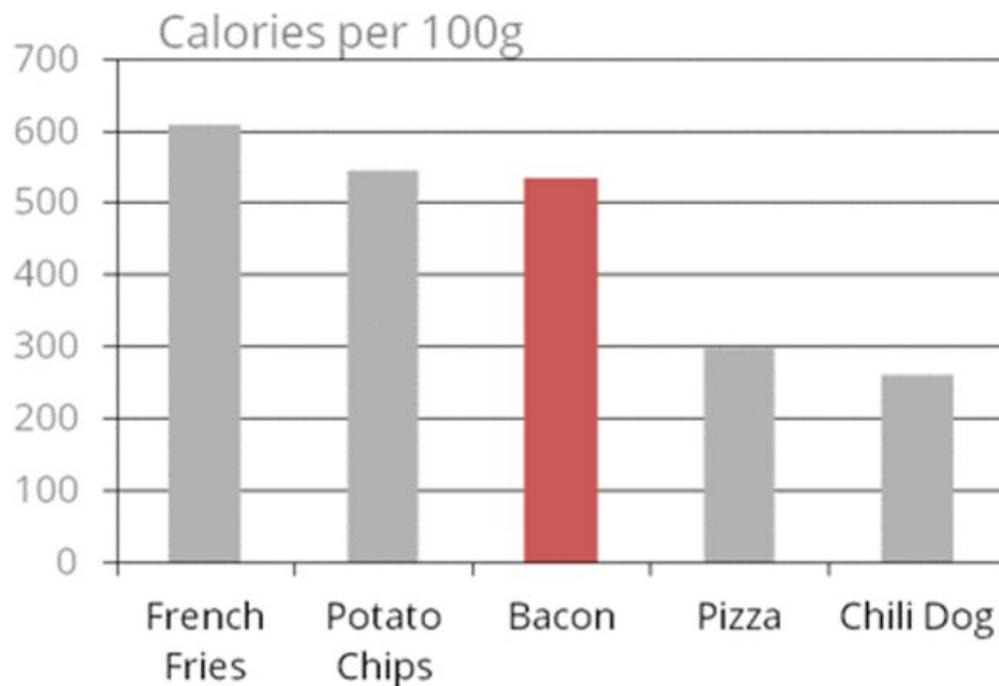
Lighten labels



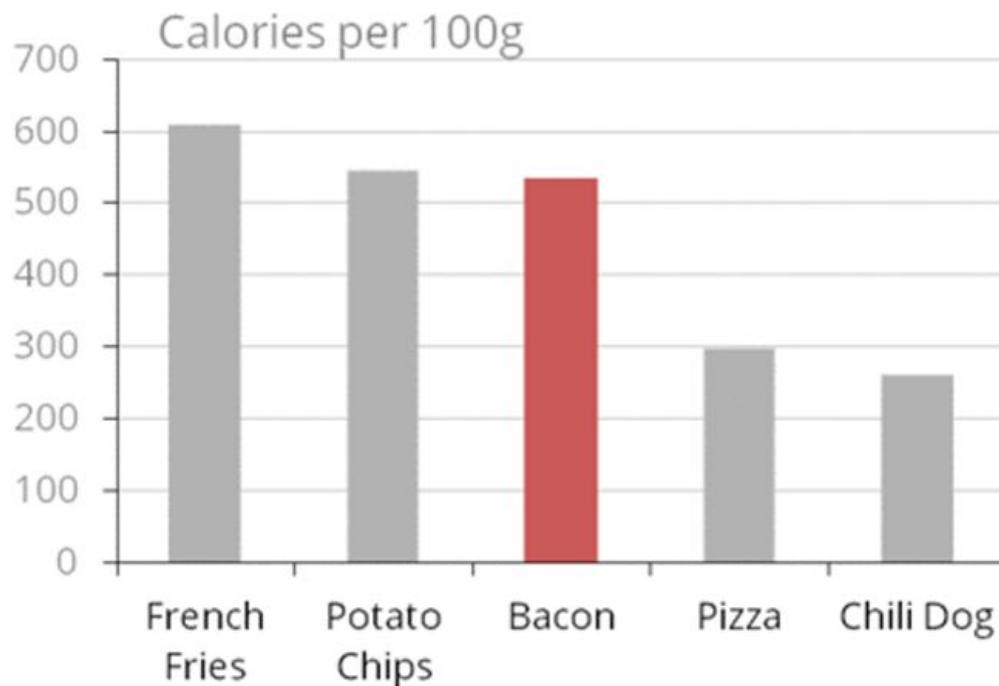
Lighten labels



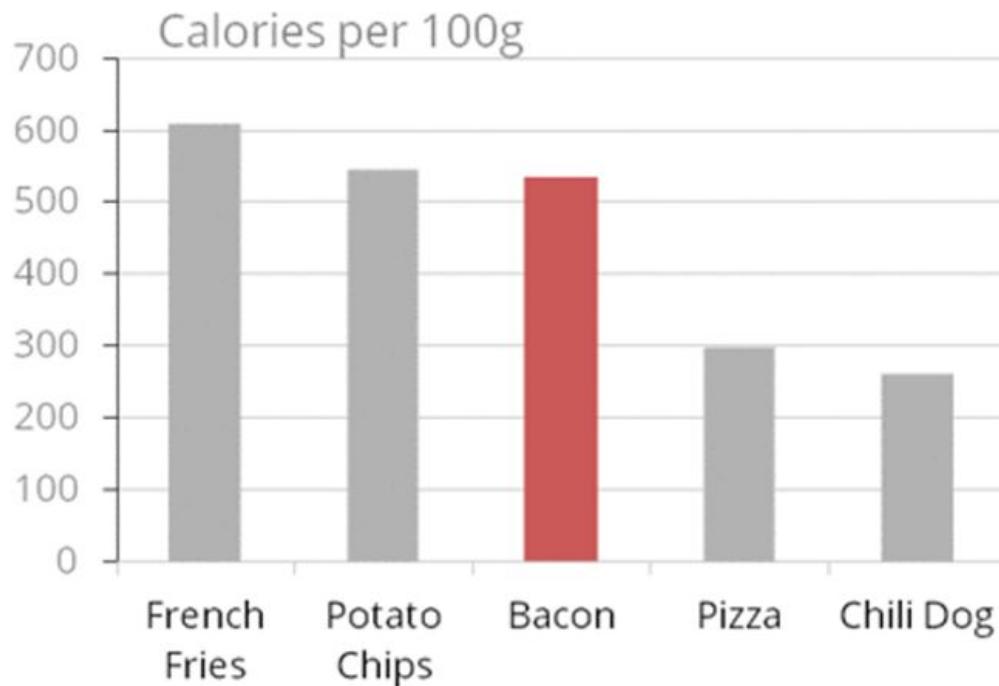
Lighten lines



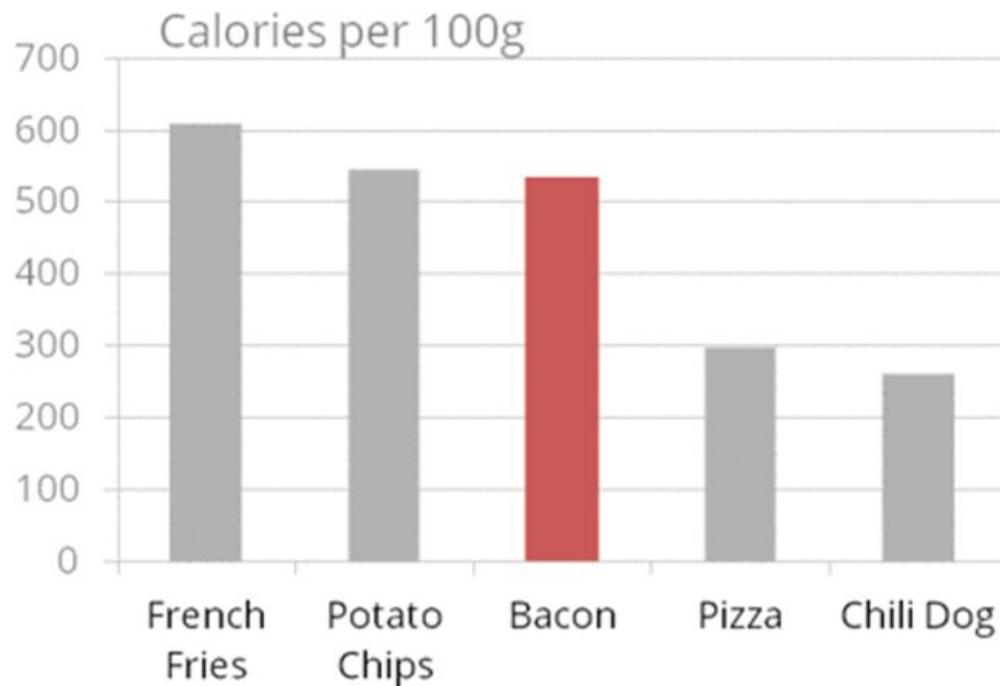
Lighten lines



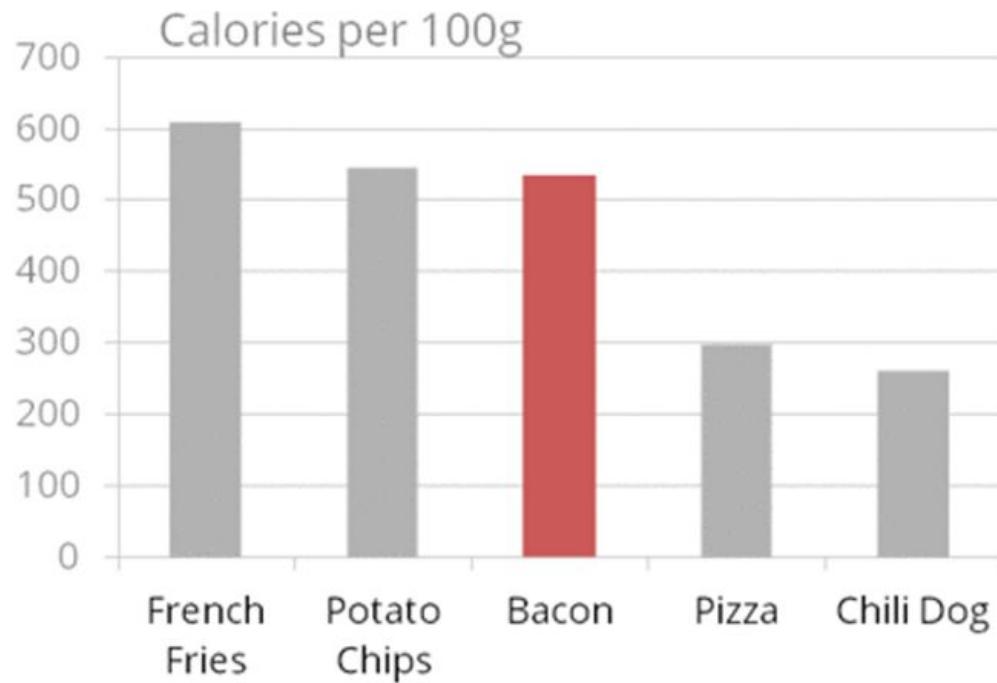
Lighten lines



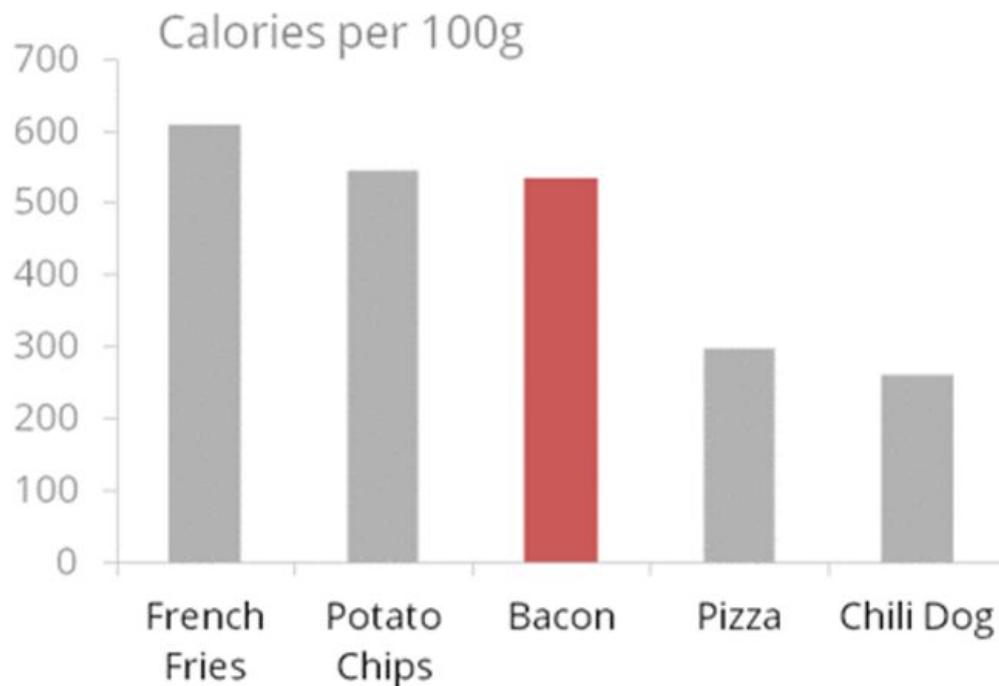
Lighten lines



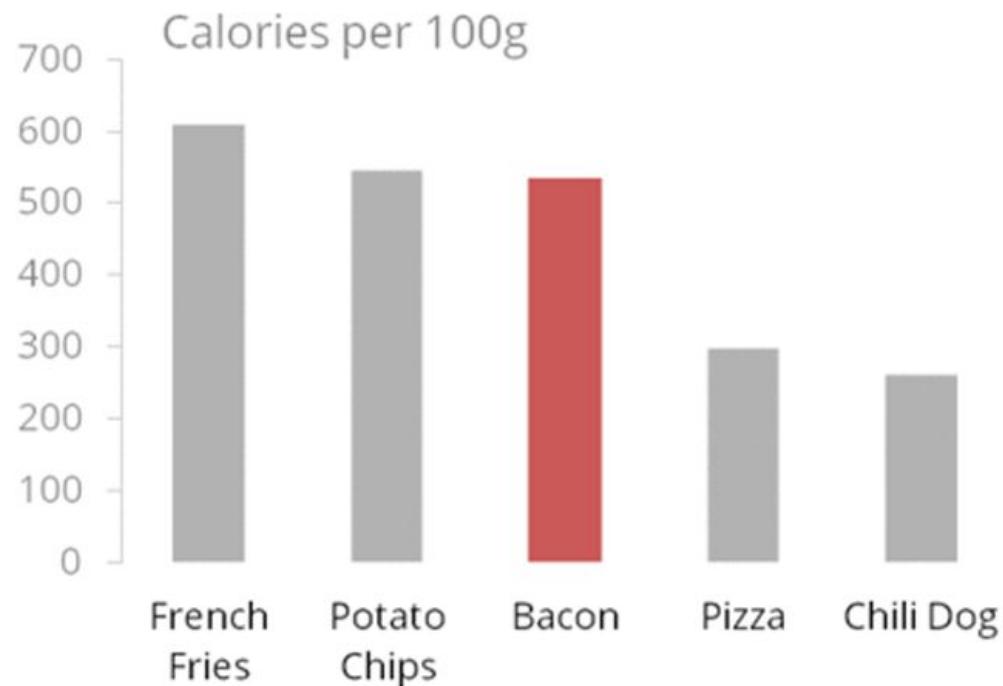
Or remove lines



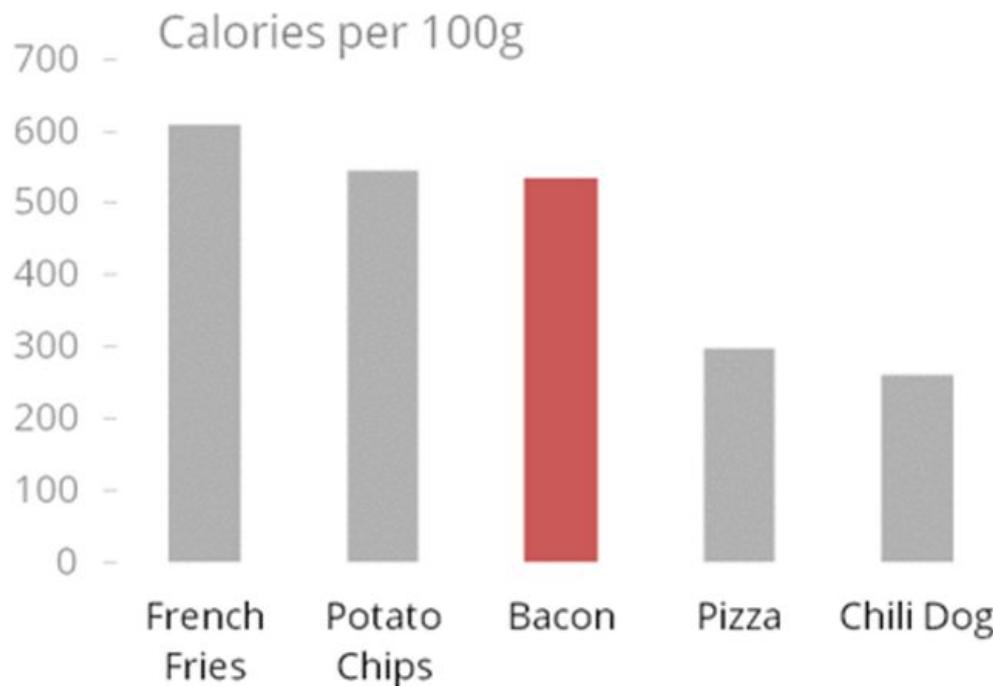
Or remove lines



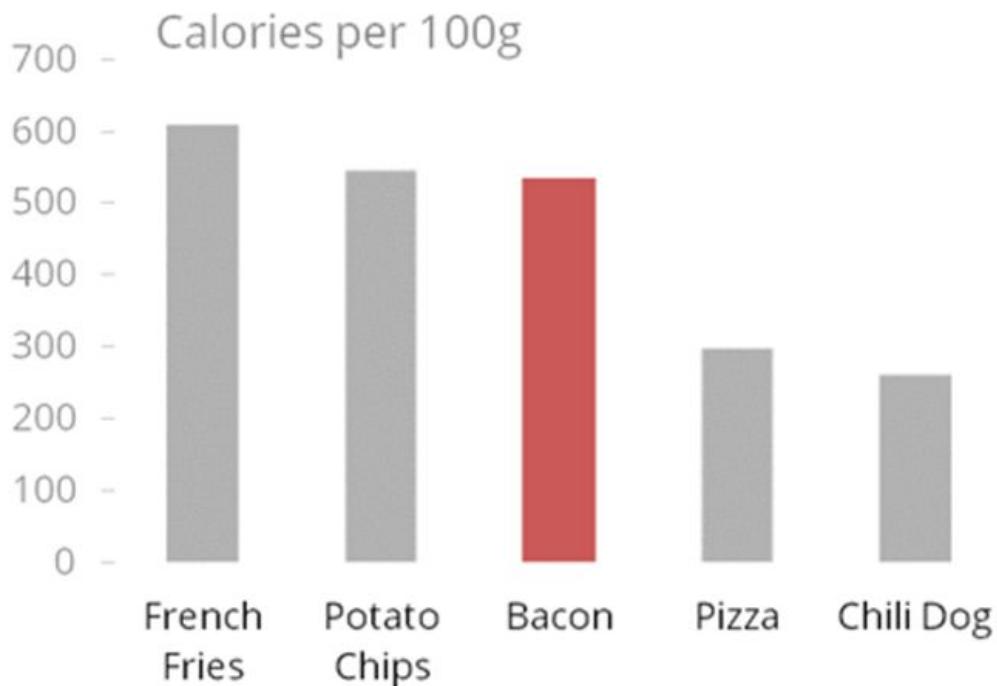
Or remove lines



Or remove lines

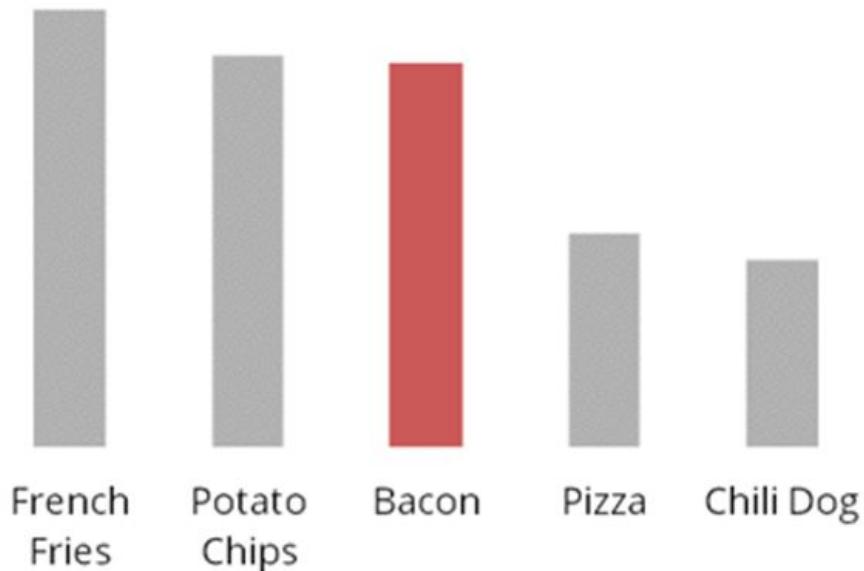


Direct label



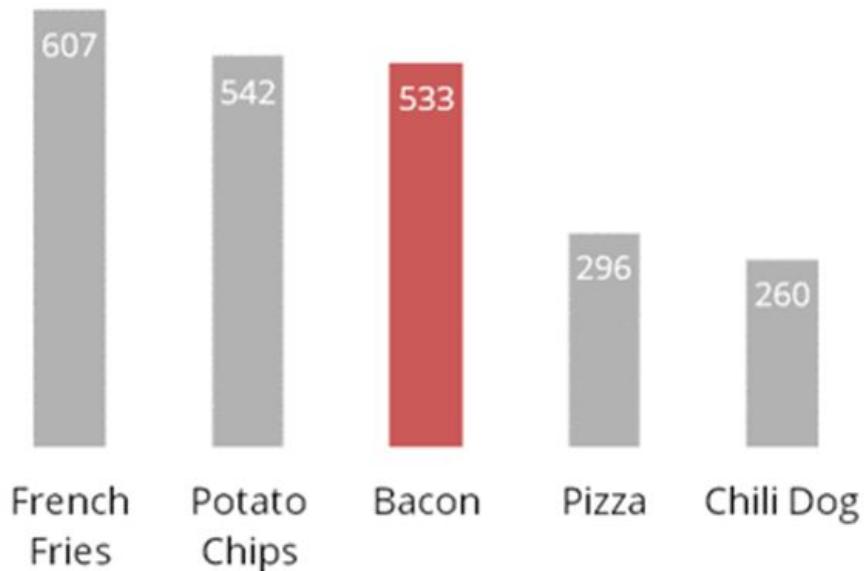
Direct label

Calories per 100g

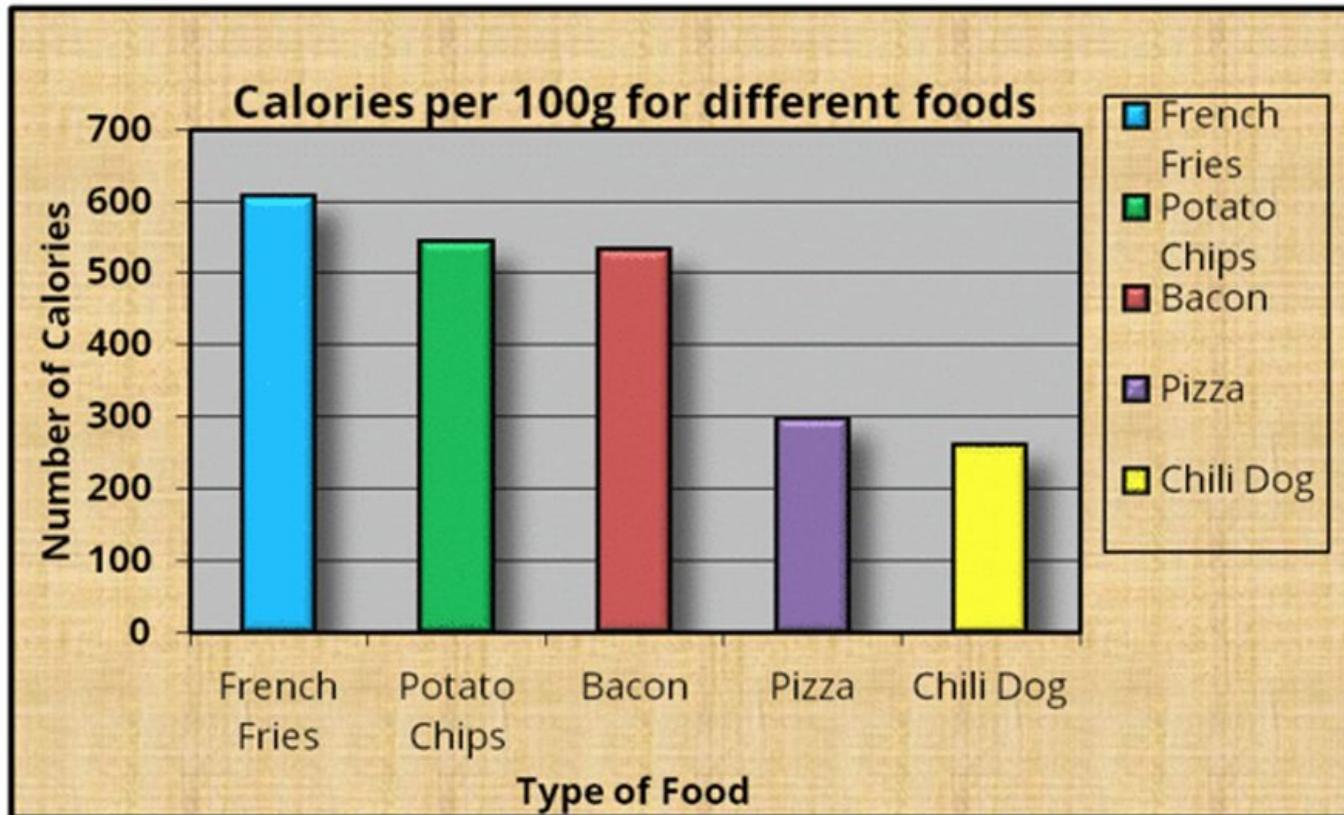


Direct label

Calories per 100g

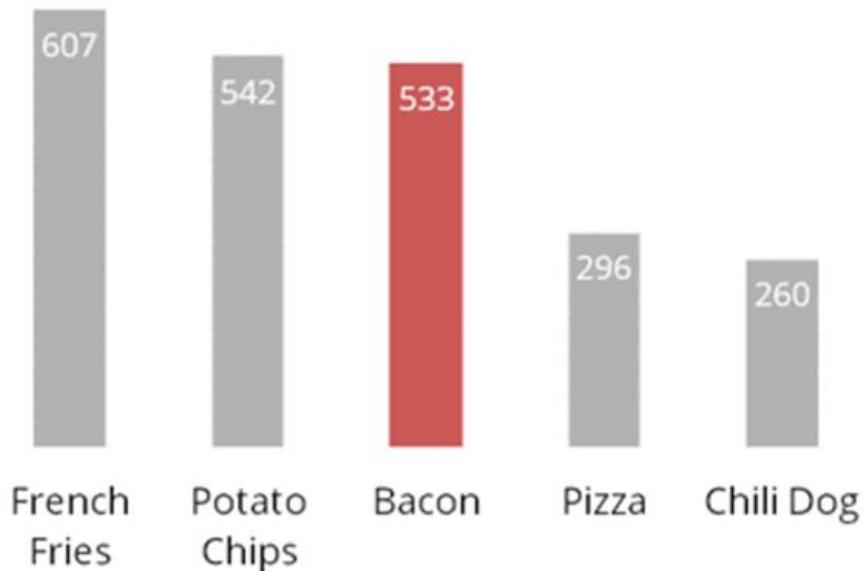


Before



After

Calories per 100g



less
is more

(effective)

(attractive)

(impactive)

adapted from Brad Voytek