



DATA SCIENCE

¿Cómo te engañan con datos?

En DS pasamos de tener **datos** a obtener **información**.

En este taller vamos a ver:

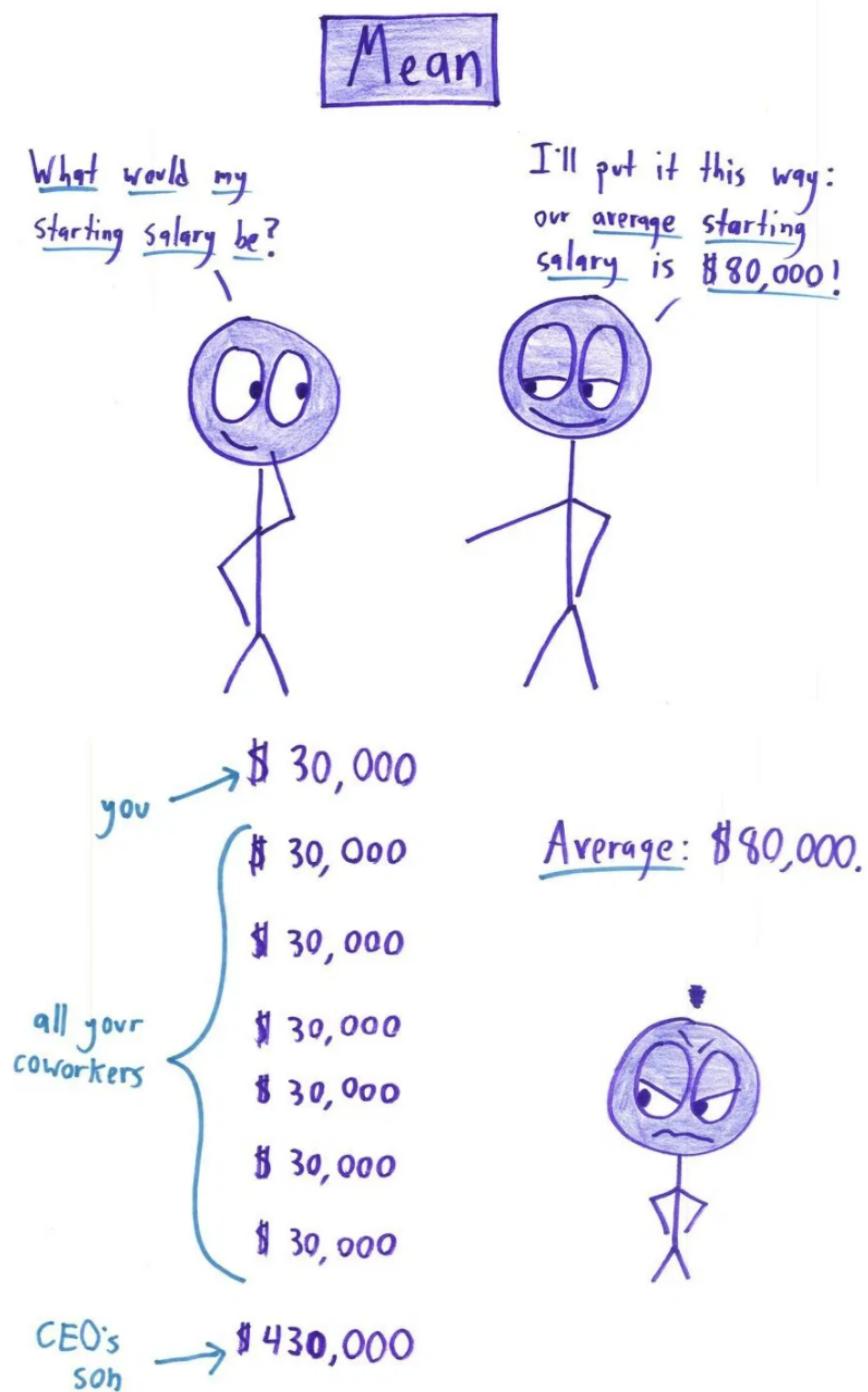
- Buenos y malos estadísticos para resumir los datos
- Buenas y malas prácticas de visualización
- Paradojas matemáticas de datos
- Equidad aplicada a modelos

A veces juegan con nuestros sesgos en la percepción para presentarnos los datos de una forma que nos va a inducir a interpretarlos erróneamente o de manera poco parcial.

EMOSIDO
ENGañADO

Estadísticos descriptivos

Veamos cómo nos pueden engañar con la estadística descriptiva, ¡empezando por la famosa media!

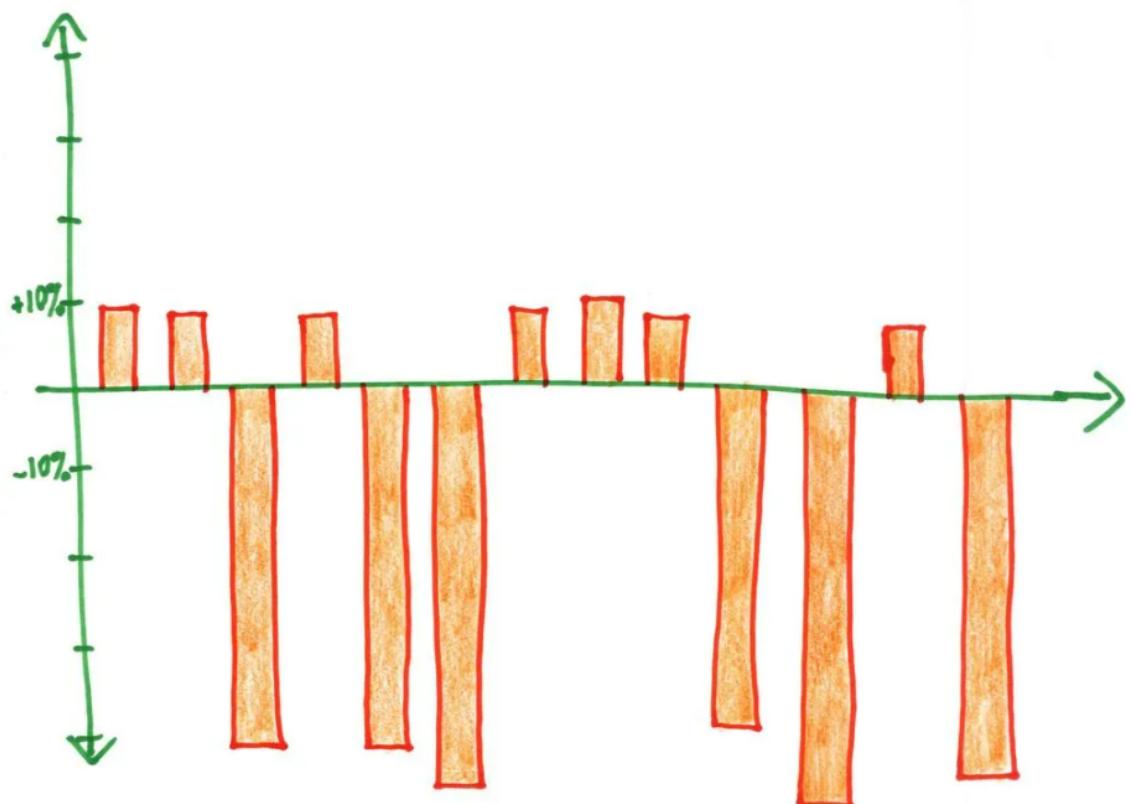


Median

So, why should I invest with you?

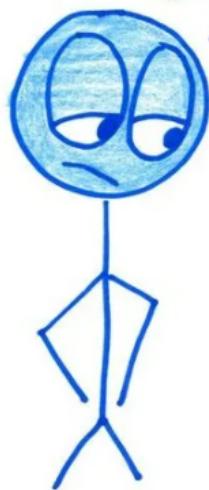


Well, not to brag, but my fund has a median gain of 8% per year!



Mode

How are you doing
on your tests?



My modal category
is 70-80%!



Score Category	Number of Tests
90s	0
80s	0
70s	2
60s	1
50s	1
40s	1
30s	1
20s	1

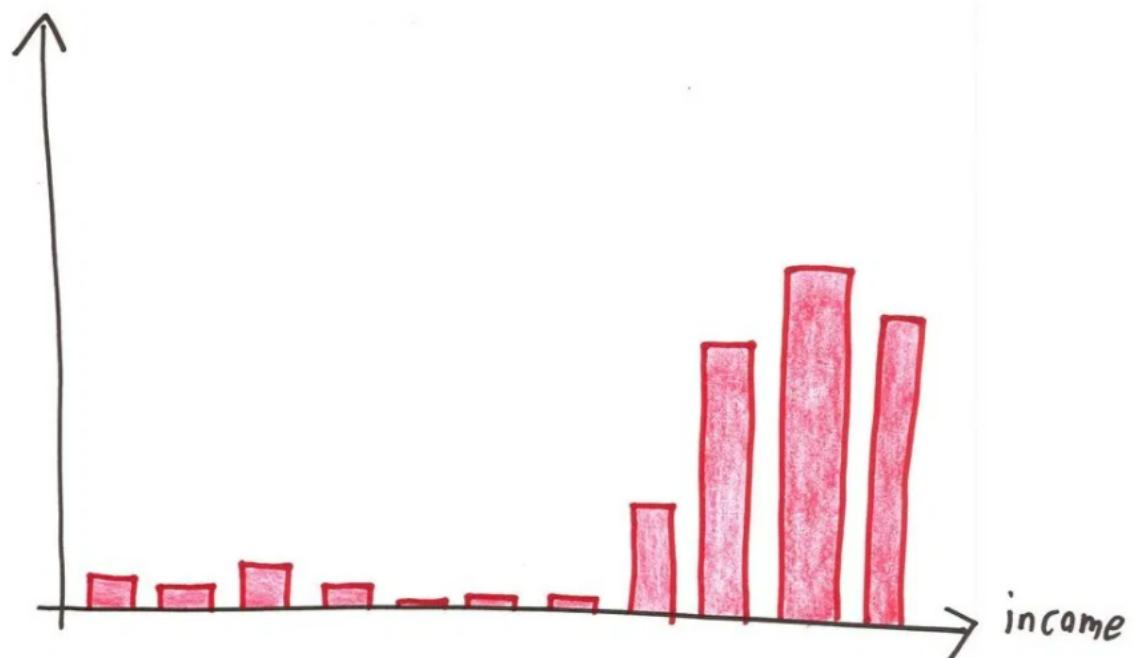


Range

Our students come from a
wide range of
Socioeconomic
backgrounds...

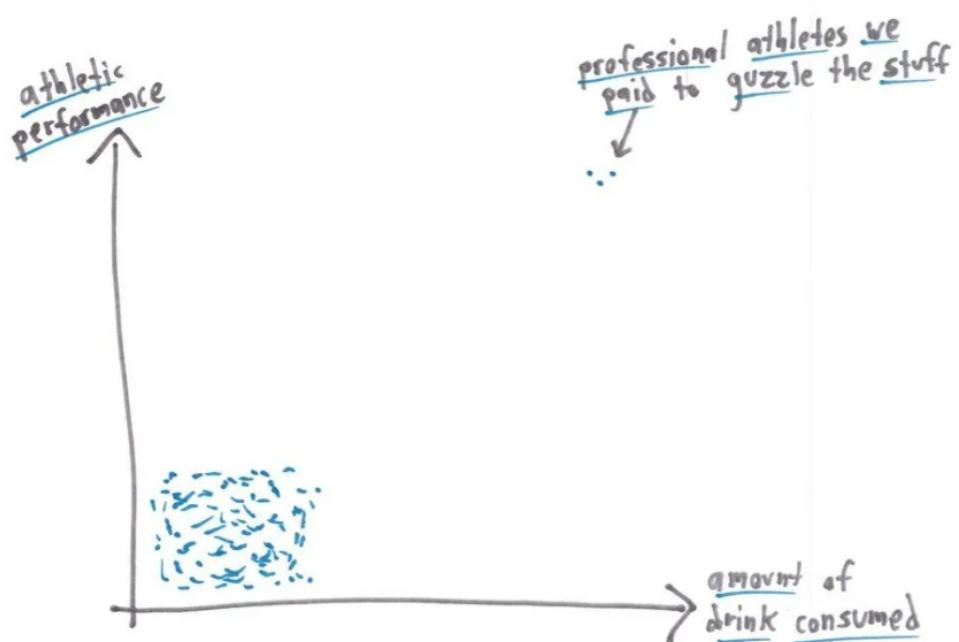
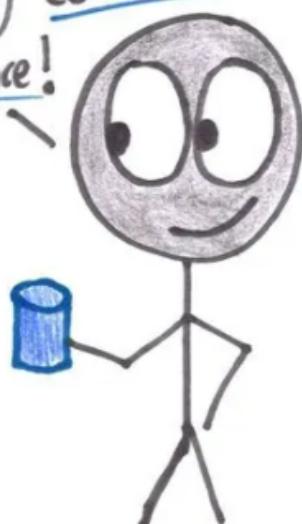


Number
of students



Correlation Coefficient

Try our energy drink —
it's highly correlated with
performance!

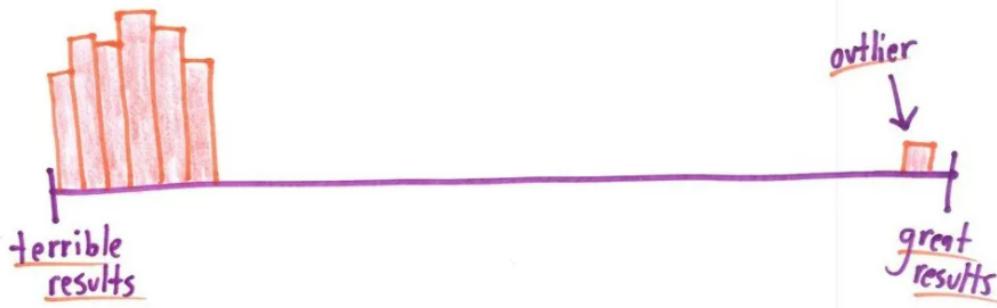


Variance

These results are a disaster!



Sure, they look bad, but there's a lot of variance!
Don't rush to judgment.



¿Y por qué SIEMPRE nos dan el dato MEDIO?

 La Voz de Galicia
La rebaja del IRPF: un euro al mes para el sueldo medio y 40 ...



... su base liquidable general fuera de 20.000 euros anuales (prácticamente la correspondiente al sueldo medio de los gallegos en el 2020).

Hace 1 día

 Libre Mercado
El sueldo medio en el sector público ya supera al privado en 800 euros mensuales



El sueldo medio en el sector público ya supera al privado en 800 euros mensuales. Institución Futuro denuncia que las subidas del SMI afectan...

Hace 3 semanas

 La Vanguardia
Alarma sindical por los bajos salarios del talento tecnológico
En cuanto a sueldos, no hay estadística oficial. ... Según la asociación Barcelona Digital Talent, el sueldo medio de un profesional digital...



Hace 4 días

 Vozpópolis
Los sueldos no acompañan a la inflación: son un 1,1% inferiores a los de hace 20 años
Si los salarios hubieran subido lo mismo que los precios, el sueldo medio anual debería situarse en los 27.773 euros. Esto supone que los...



Hace 3 semanas

 elEconomista.es
Este es el salario medio de los profesores por países: España no aparece entre los 10 primeros
Este es el salario medio de los profesores por países: España no aparece entre los 10 primeros. Luxemburgo, donde los educadores ganan 101.360...

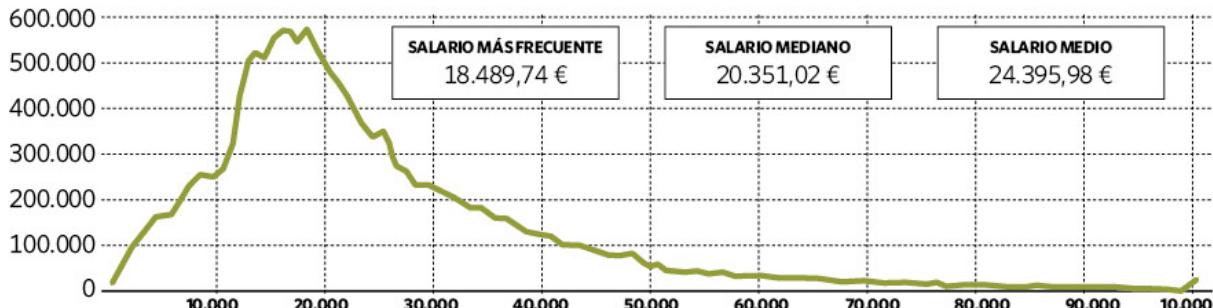


Hace 3 semanas

La media no te dice **NADA** de la dispersión de los datos.

Distribución Ganancia bruta 2019

A la izquierda, número de trabajadores y abajo, euros

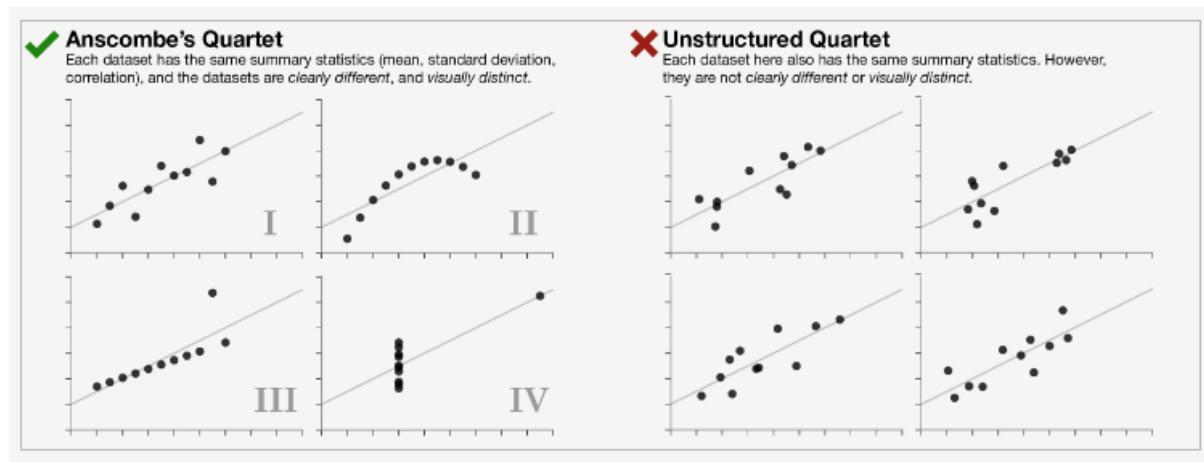


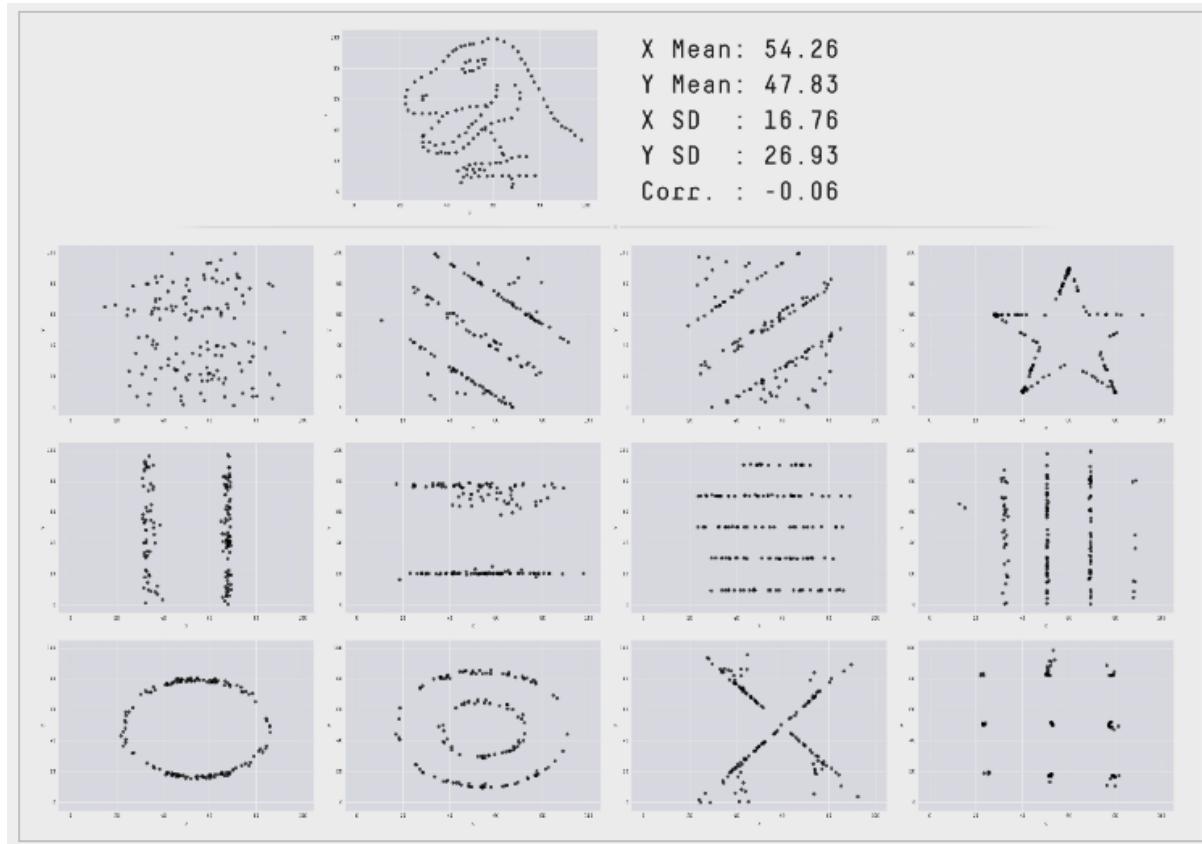
Fuente: Encuesta Anual de Estructura Salarial. Año 2019.

elEconomista

¿Y noticias con el sueldo mediano, para cuándo?

Entonces, si conozco la media y la varianza, ya sé mucho sobre mis datos, ¿no?





También juegan con nuestros **sesgos mentales**:

¿Cuánta gente distinta crees que tendrías que juntar aleatoriamente para que dos de esas personas cumplan años el mismo día?

Solo con 23 personas ya tienes un 50% de probabilidad de tener 2 personas con el mismo cumpleaños. Solemos pensar que necesitamos muchas más personas, porque SOLO pensamos en nosotros mismos y no en la combinación de cada persona con los demás.

¿Y si vas a un concurso de televisión?

The Monty Hall Problem



¿Te quedas con tu puerta o la cambias?

Te están engañando si no conoces la probabilidad CONDICIONADA.

La probabilidad de acertar a la primera es $\frac{1}{3}$ y NUNCA CAMBIA.

Tras abrir una puerta, tu intento inicial SIGUE SIENDO $\frac{1}{3}$, luego la puerta cerrada tendrá una probabilidad de $\frac{2}{3}$

¿Quieres verlo más claro? Piensa en 1000 puertas iniciales con un coche. Tras elegir una puerta, te abren todas las restantes menos una. Tienes una probabilidad de acertar a la primera de $1/1000 = 0.001$ y una probabilidad de que esté en la otra puerta cerrada de 0.999. Ahora lo vemos más claro.

Buenas y malas prácticas en visualización

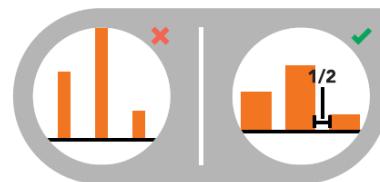
BAR CHART

DESIGN BEST PRACTICES



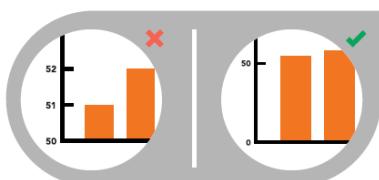
USE HORIZONTAL LABELS

Avoid steep diagonal or vertical type, as it can be difficult to read.



SPACE BARS APPROPRIATELY

Space between bars should be $\frac{1}{2}$ bar width.



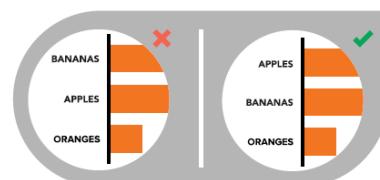
START THE Y-AXIS VALUE AT 0

Starting at a value above zero truncates the bars and doesn't accurately reflect the full value.



USE CONSISTENT COLORS

Use one color for bar charts. You may use an accent color to highlight a significant data point.



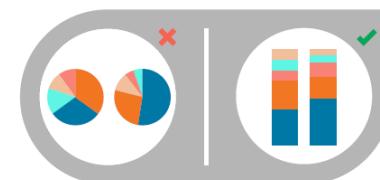
ORDER DATA APPROPRIATELY

Order categories alphabetically, sequentially, or by value.



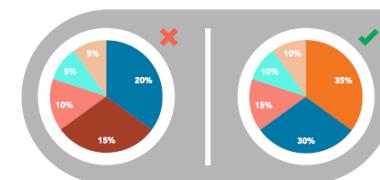
VISUALIZE NO MORE THAN 5 CATEGORIES PER CHART

It is difficult to differentiate between small values; depicting too many slices decreases the impact of the visualization. If needed, you can group smaller values into an "other" or "miscellaneous" category, but make sure it does not hide interesting or significant information.



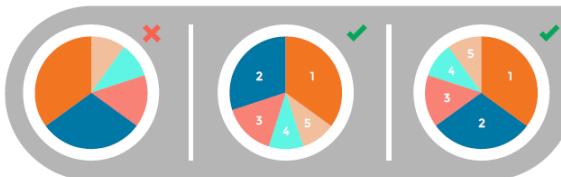
DON'T USE MULTIPLE PIE CHARTS FOR COMPARISON

Slice sizes are very difficult to compare side-by-side. Use a stacked bar chart instead.



MAKE SURE ALL DATA ADDS UP TO 100%

Verify that values total 100% and that pie slices are sized proportionate to their corresponding value.



ORDER SLICES CORRECTLY

There are two ways to order sections, both of which are meant to aid comprehension:

OPTION 1
Place the largest section at 12 o'clock, going clockwise. Place the second largest section at 12 o'clock, going counterclockwise. The remaining sections can be placed below, continuing counterclockwise.

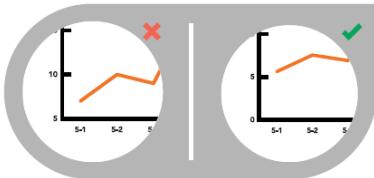
OPTION 2
Start the largest section at 12 o'clock, going clockwise. Place remaining sections in descending order, going clockwise.

PIE CHART

DESIGN BEST PRACTICES

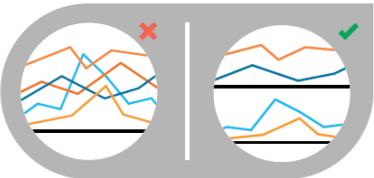
LINE CHART

DESIGN BEST PRACTICES



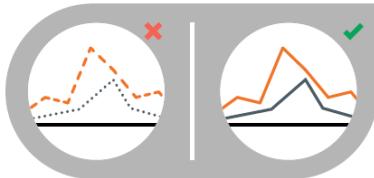
INCLUDE A ZERO BASELINE IF POSSIBLE

Although a line chart does not have to start at a zero baseline, it should be included if possible. If relatively small fluctuations in data are meaningful (e.g., in stock market data), you may truncate the scale to showcase these variances.



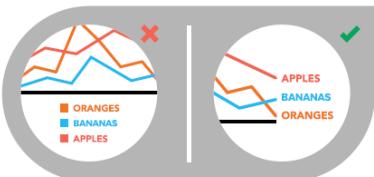
DON'T PLOT MORE THAN 4 LINES

If you need to display more, break them out into separate charts for better comparison.



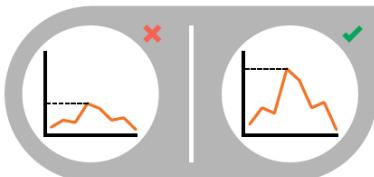
USE SOLID LINES ONLY

Dashed and dotted lines can be distracting.



LABEL THE LINES DIRECTLY

This lets readers quickly identify lines and corresponding labels instead of referencing a legend.

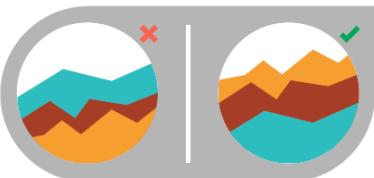


USE THE RIGHT HEIGHT

Plot all data points so that the line chart takes up approximately two-thirds of the y-axis' total scale.

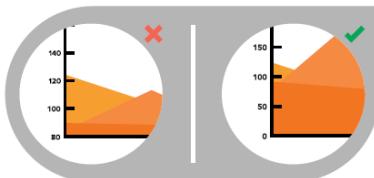
AREA CHART

DESIGN BEST PRACTICES



MAKE IT EASY TO READ

In stacked area charts, arrange data to position categories with highly variable data on the top of the chart and low variability on the bottom.



START Y-AXIS VALUE AT 0

Starting the axis above zero truncates the visualization of values.



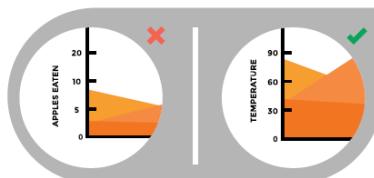
DON'T DISPLAY MORE THAN 4 DATA CATEGORIES

Too many will result in a cluttered visual that is difficult to decipher.



USE TRANSPARENT COLORS

In standard area charts, ensure data isn't obscured in the background by ordering thoughtfully and using transparency.

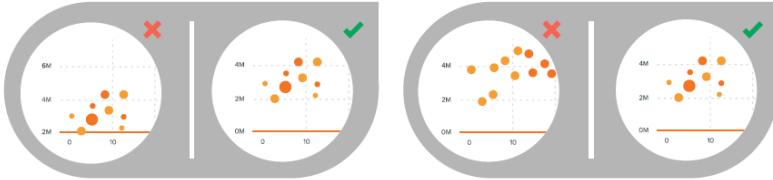


DON'T USE AREA CHARTS TO DISPLAY DISCRETE DATA

The connected lines imply intermediate values, which only exist with continuous data.

SCATTER PLOT

DESIGN BEST PRACTICES

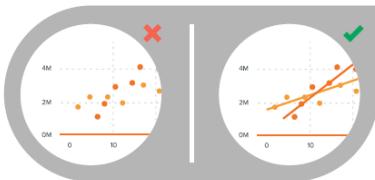


START Y-AXIS VALUE AT 0

Starting the axis above zero truncates the visualization of values.

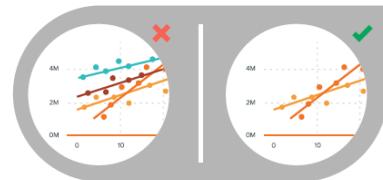
INCLUDE MORE VARIABLES

Use size and dot color to encode additional data variables.



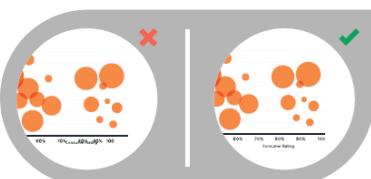
USE TREND LINES

These help draw correlation between the variables to show trends.



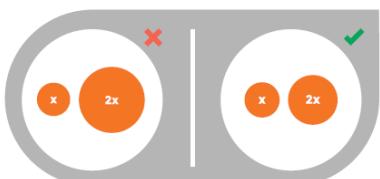
DON'T COMPARE MORE THAN 2 TREND LINES

Too many lines make data difficult to interpret.



MAKE SURE LABELS ARE VISIBLE

All labels should be unobstructed and easily identified with the corresponding bubble.



SIZE BUBBLES APPROPRIATELY

Bubbles should be scaled according to area, not diameter.



DON'T USE ODD SHAPES

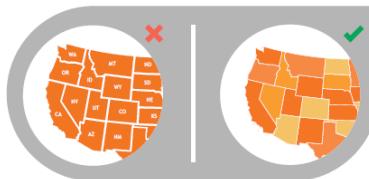
Avoid adding too much detail or using shapes that are not entirely circular; this can lead to inaccuracies.

BUBBLE CHART

DESIGN BEST PRACTICES

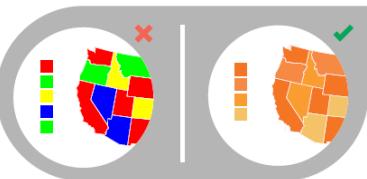
HEAT MAP

DESIGN BEST PRACTICES



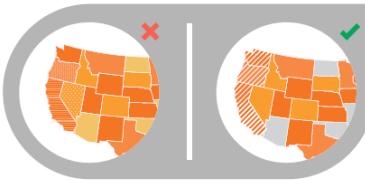
USE A SIMPLE MAP OUTLINE

These lines are meant to frame the data, not distract.



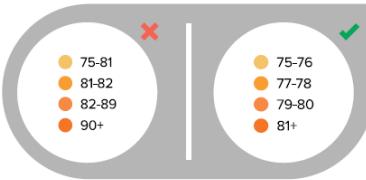
SELECT COLORS APPROPRIATELY

Some colors stand out more than others, giving unnecessary weight to that data. Instead, use a single color with varying shade or a spectrum between two analogous colors to show intensity. Also remember to intuitively code color intensity according to values.



USE PATTERNS SPARINGLY

A pattern overlay that indicates a second variable is acceptable, but using multiple is overwhelming and distracting.



CHOOSE APPROPRIATE DATA RANGES

Select 3-5 numerical ranges that enable fairly even distribution of data between them. Use +/- signs to extend high and low ranges.

10 DATA DESIGN DOS AND DON'TS

Designing your data doesn't have to be overwhelming. With a basic understanding of how different data sets should be visualized, along with a few fundamental design tips and best practices, you can create more accurate, more effective data visualizations. Follow these 10 tips to ensure your design does your data justice.



1 | DO USE ONE COLOR TO REPRESENT EACH CATEGORY.



2 | DO ORDER DATA SETS USING LOGICAL HIERARCHY.



3 | DO USE CALLOUTS TO HIGHLIGHT IMPORTANT OR INTERESTING INFORMATION.



4 | DO VISUALIZE DATA IN A WAY THAT IS EASY FOR READERS TO COMPARE VALUES.



5 | DO USE ICONS TO ENHANCE COMPREHENSION AND REDUCE UNNECESSARY LABELING.



6 | DON'T USE HIGH CONTRAST COLOR COMBINATIONS SUCH AS RED/GREEN OR BLUE/YELLOW.



7 | DON'T USE 3D CHARTS. THEY CAN SKEW PERCEPTION OF THE VISUALIZATION.



8 | DON'T ADD CHART JUNK. UNNECESSARY ILLUSTRATIONS, DROP SHADOWS, OR ORNAMENTS DISTRACT FROM THE DATA.

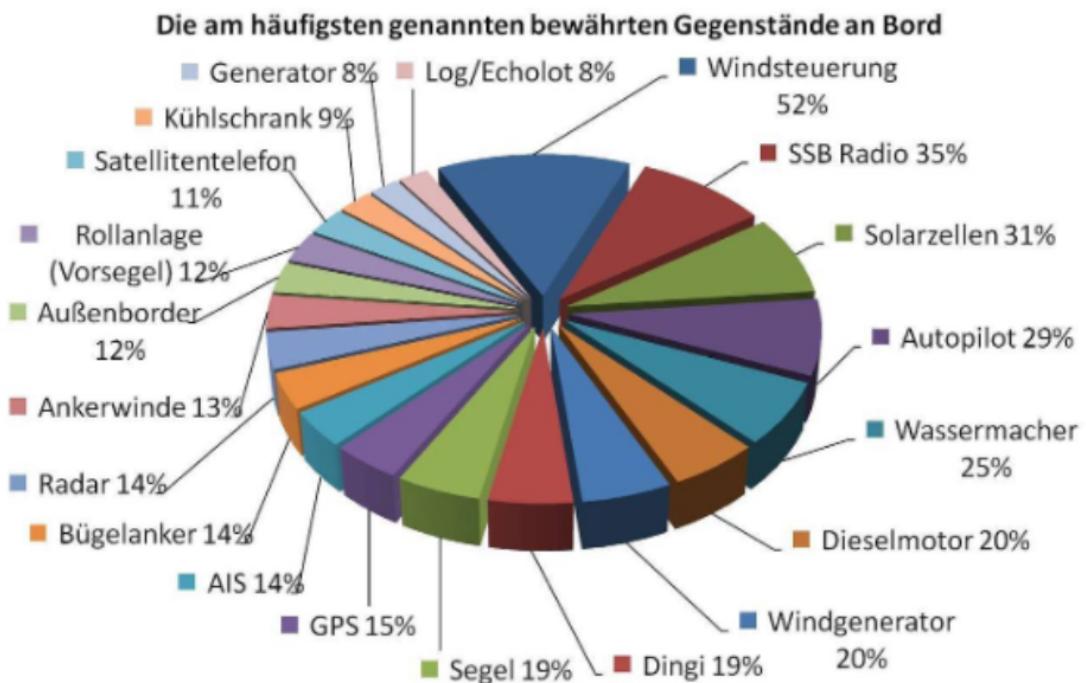


9 | DON'T USE MORE THAN 6 COLORS IN A SINGLE LAYOUT.

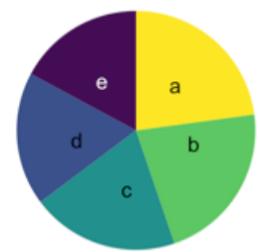
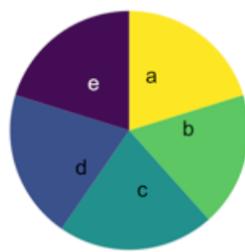
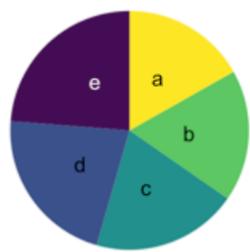


10 | DON'T USE DISTRACTING FONTS OR ELEMENTS (SUCH AS BOLD, ITALIC, OR UNDERLINED TEXT).

Quiz time! :)

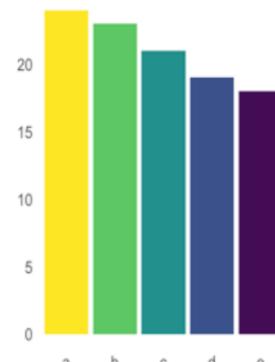
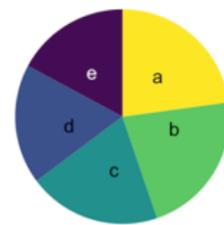
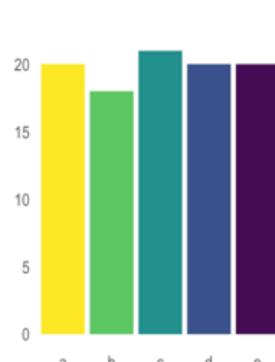
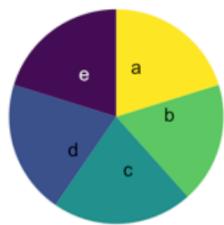
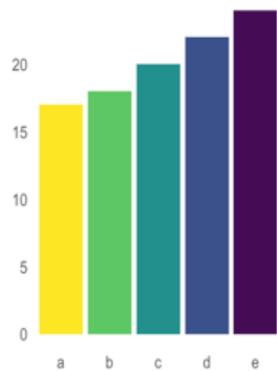
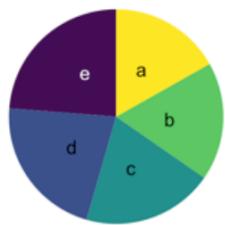


Cálculo + 3D + pie chart

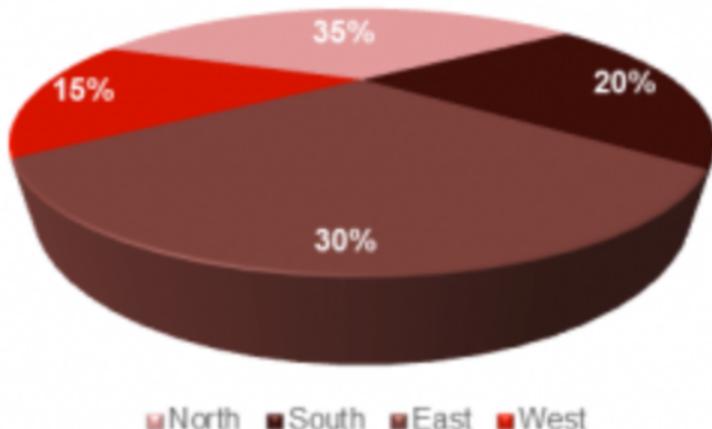


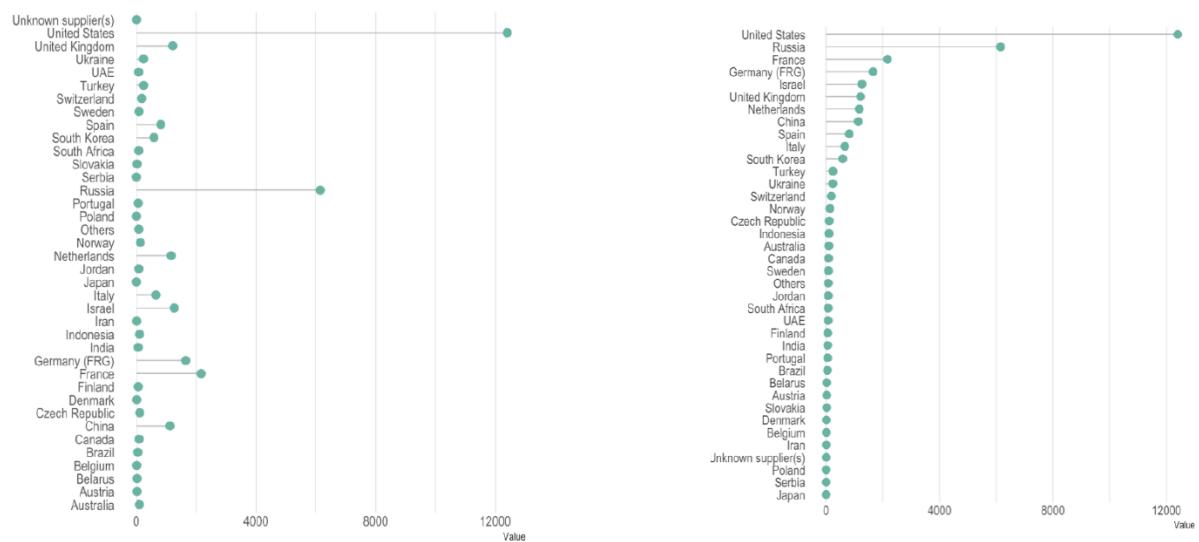
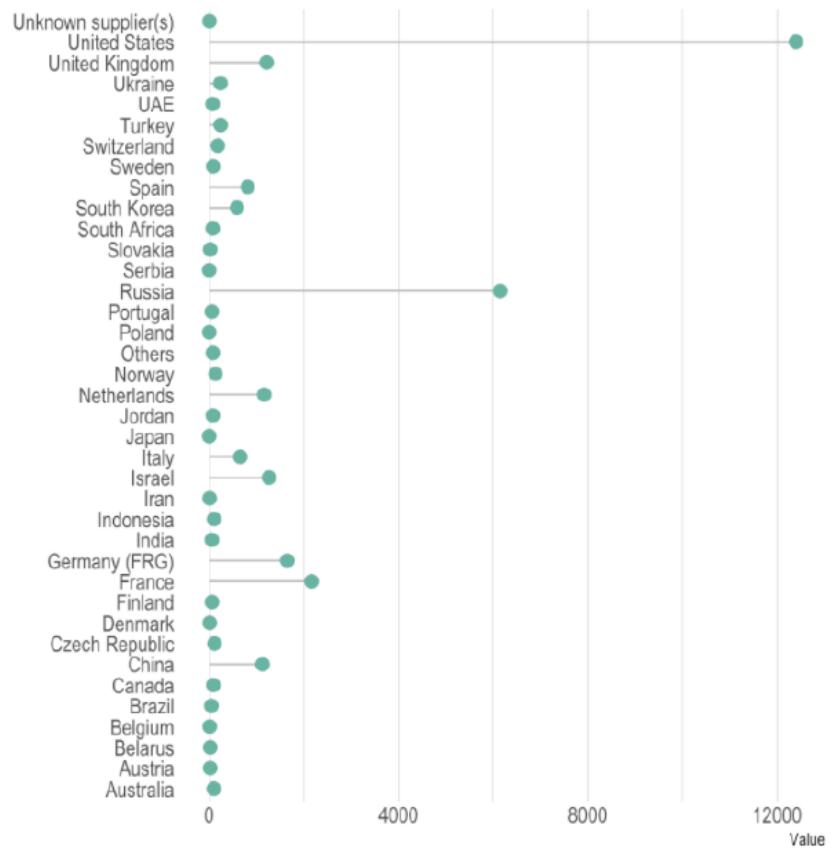
What's wrong with
pie chart?

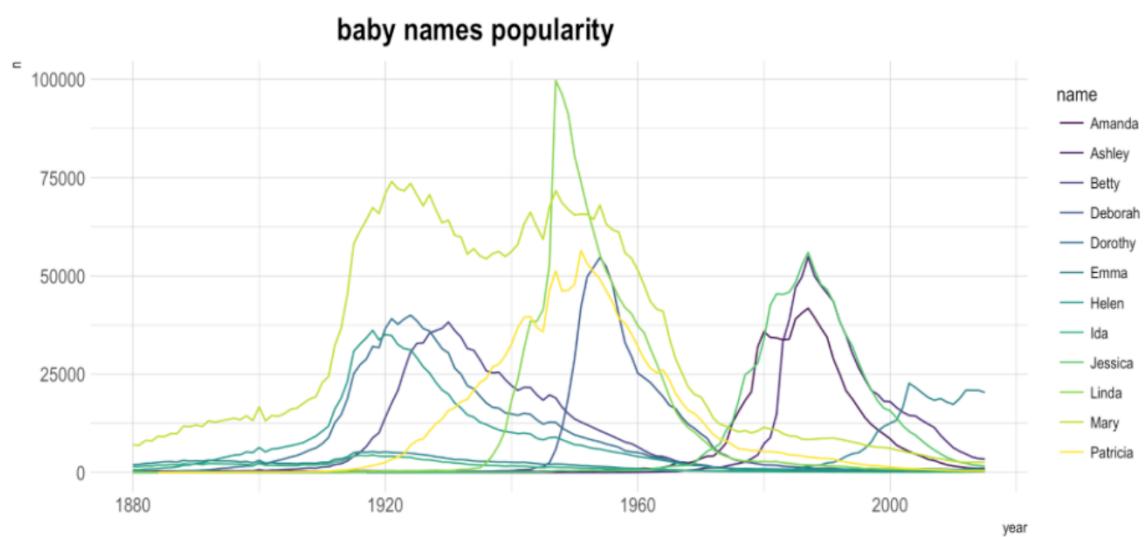
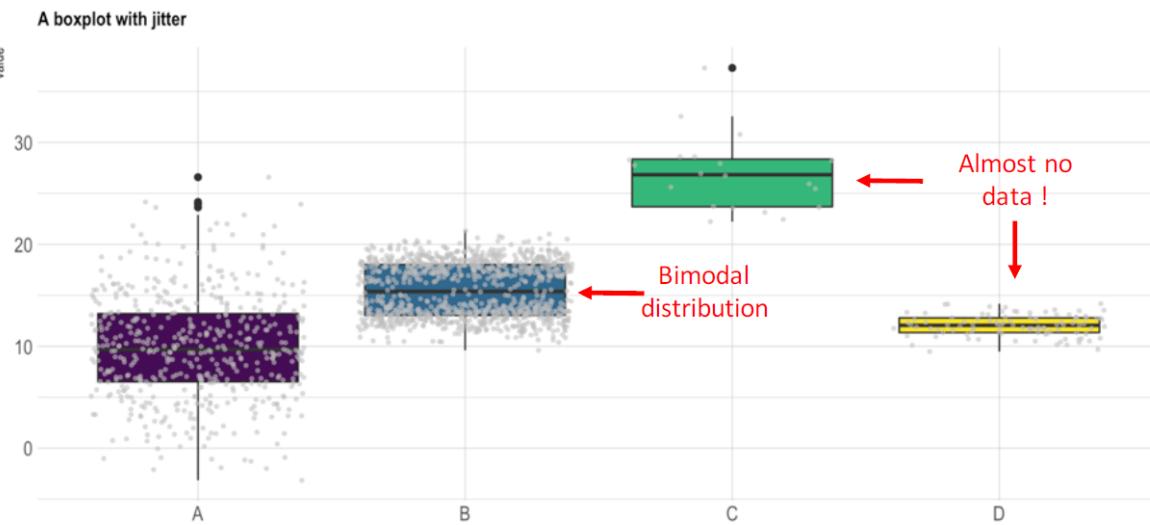
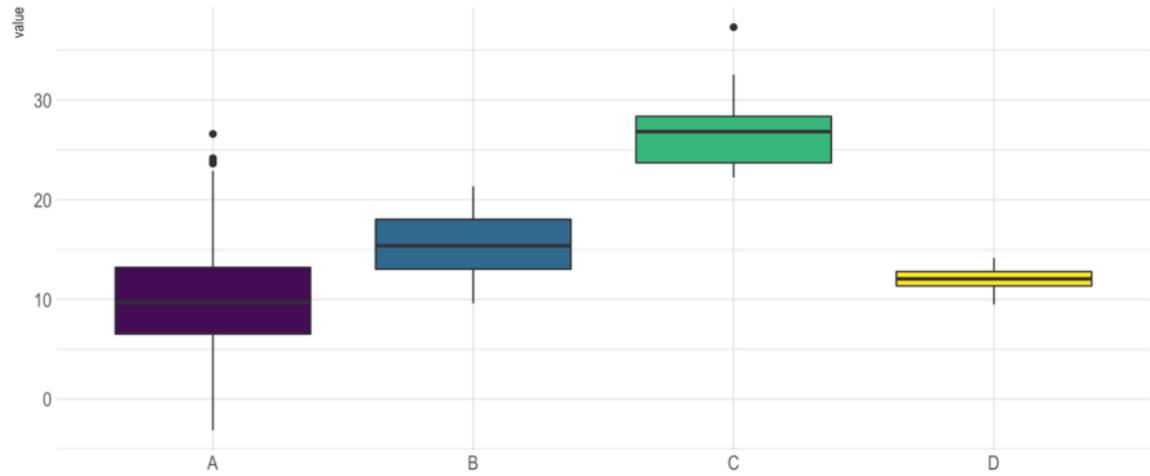
It is hard to
distinguish angles

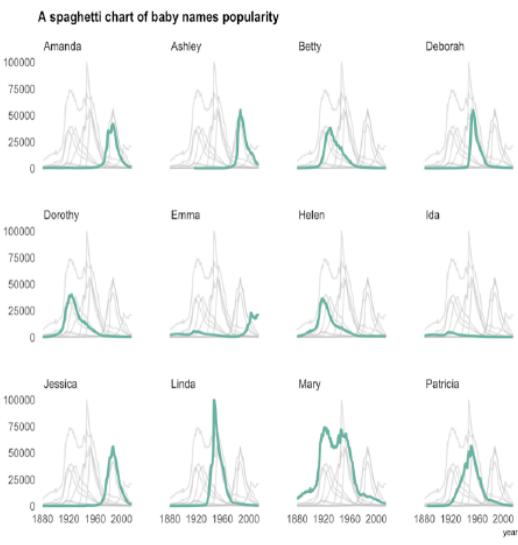
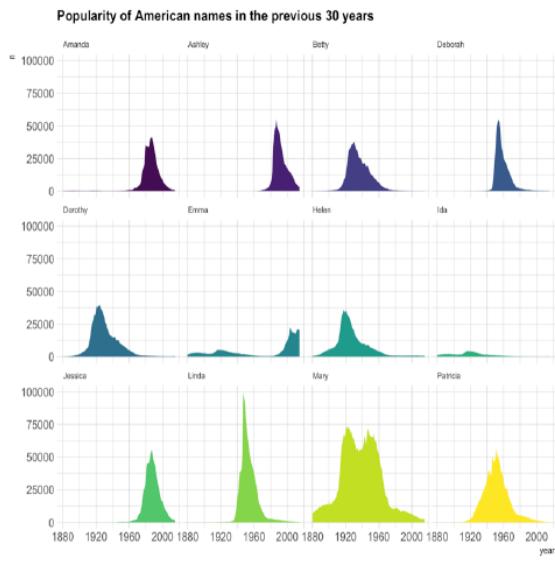
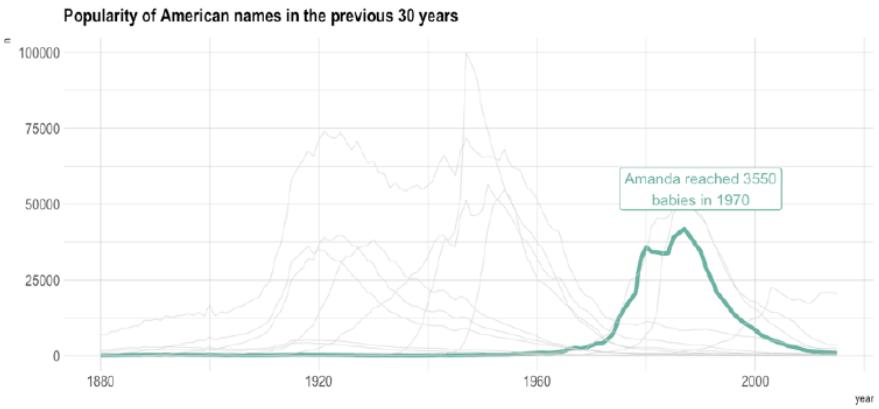


Sales Proportions By Region

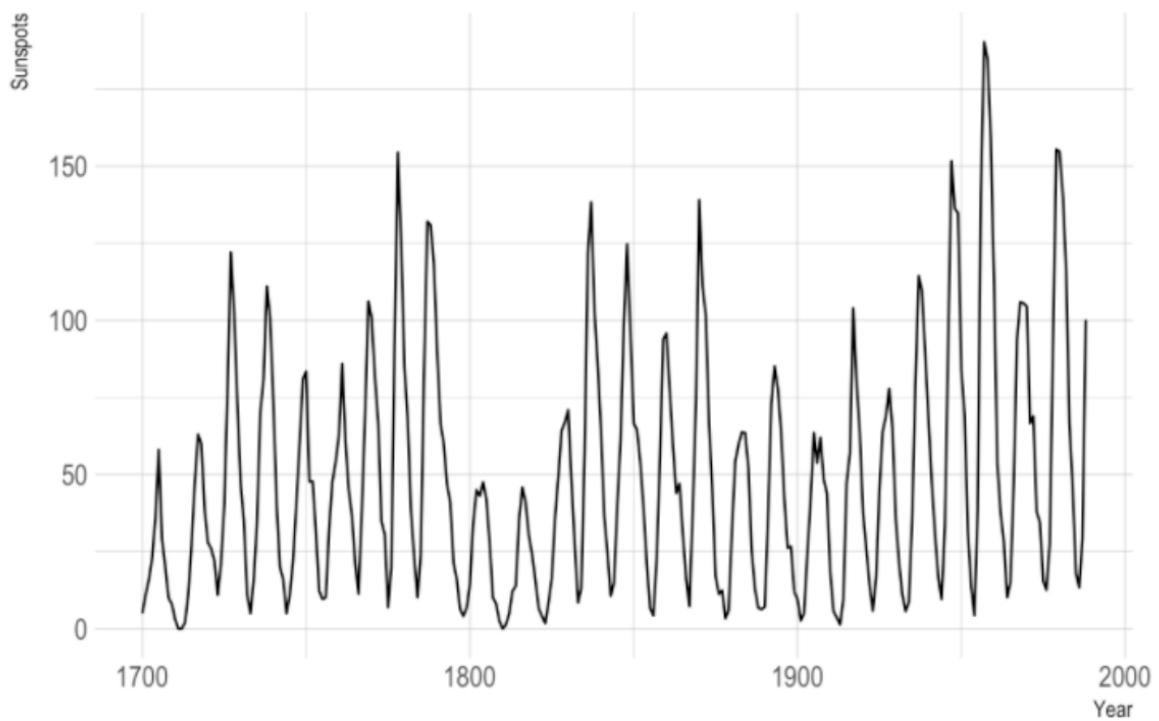




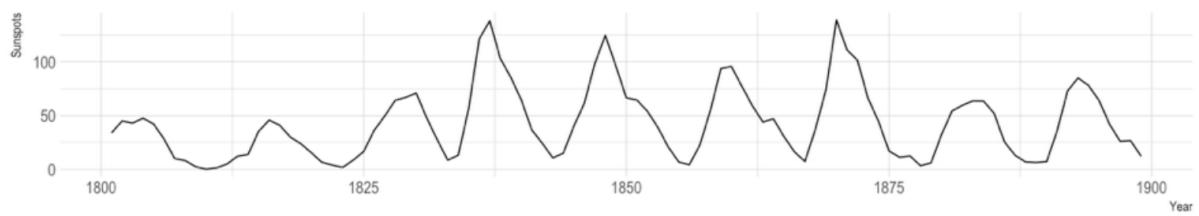




Number of sunspots per year

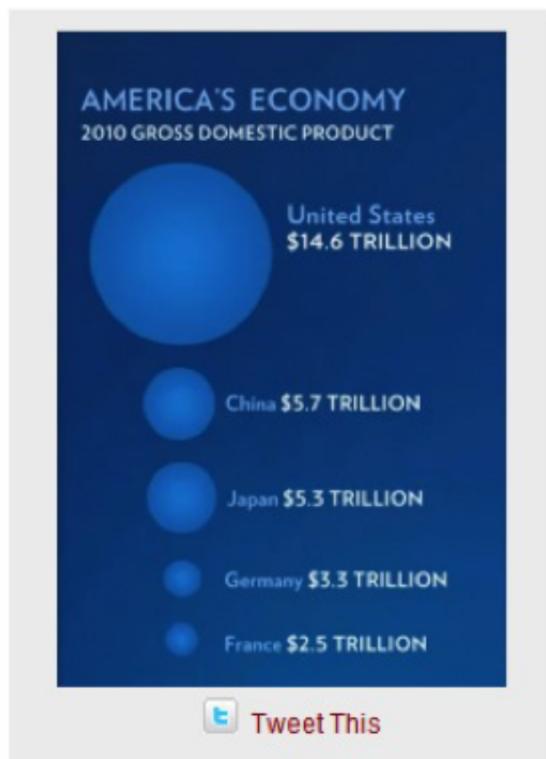


Number of sunspots per year

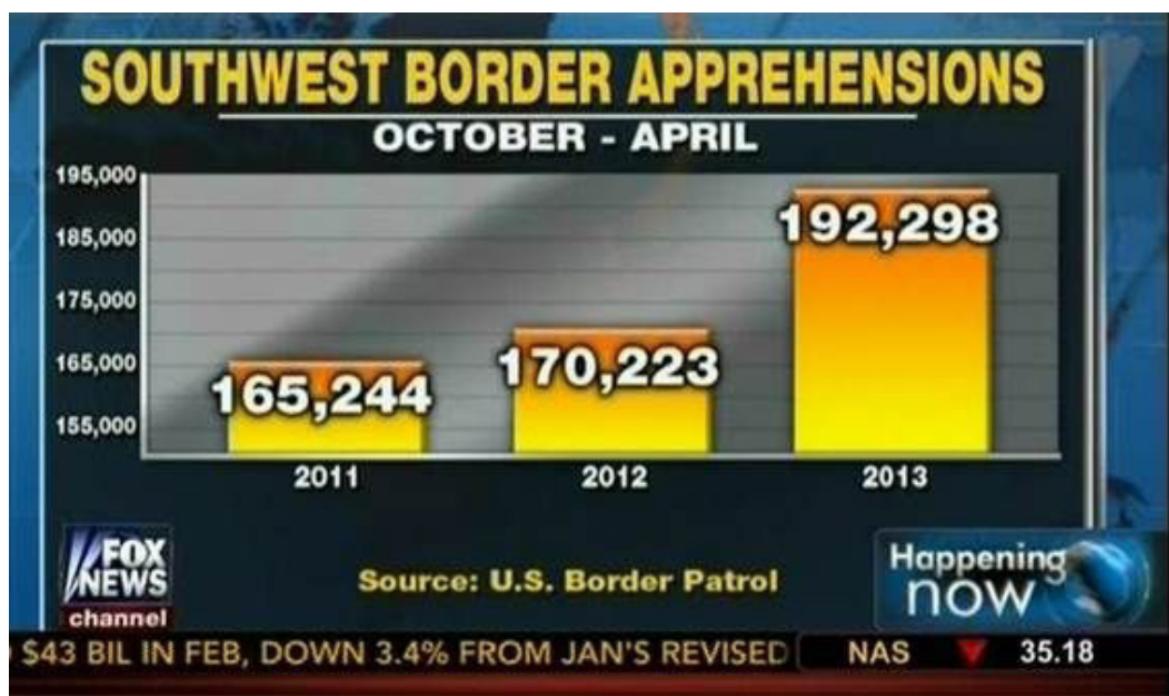


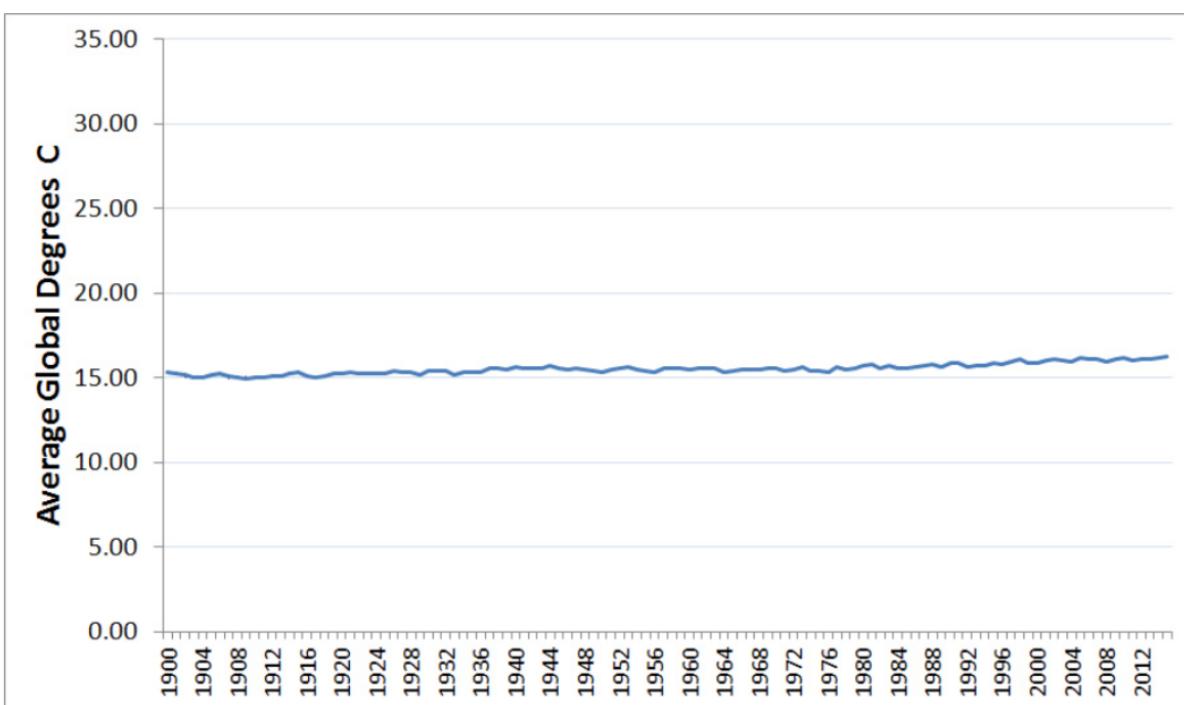
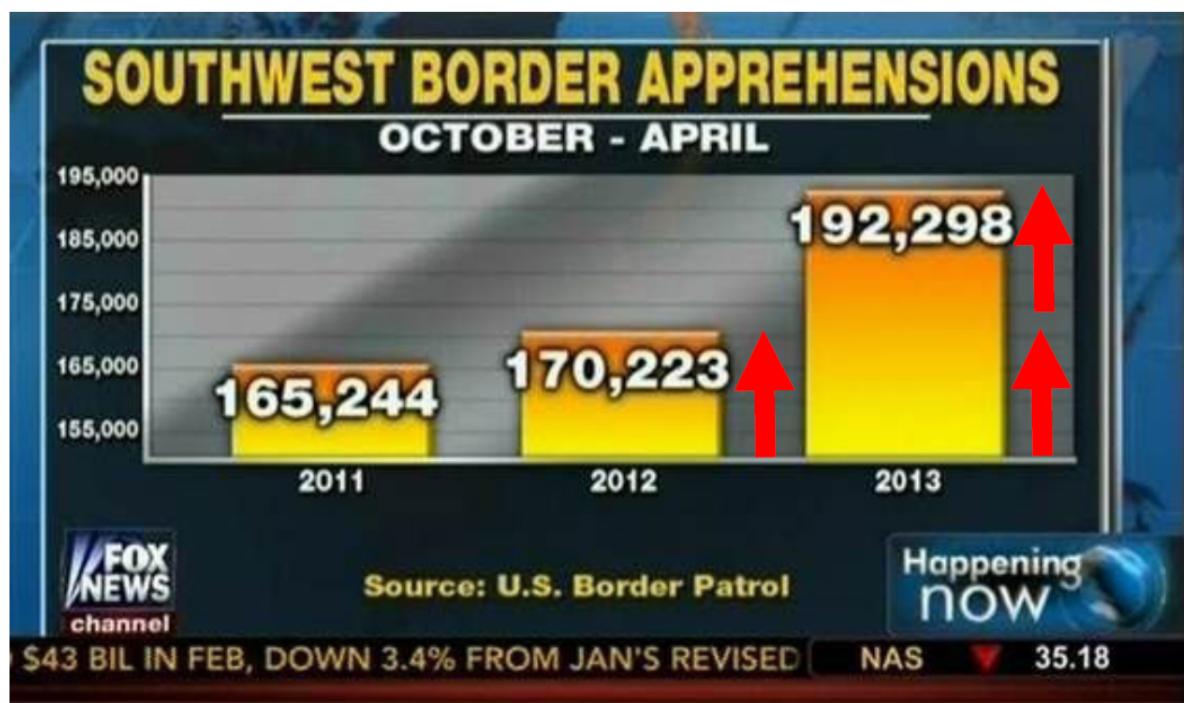


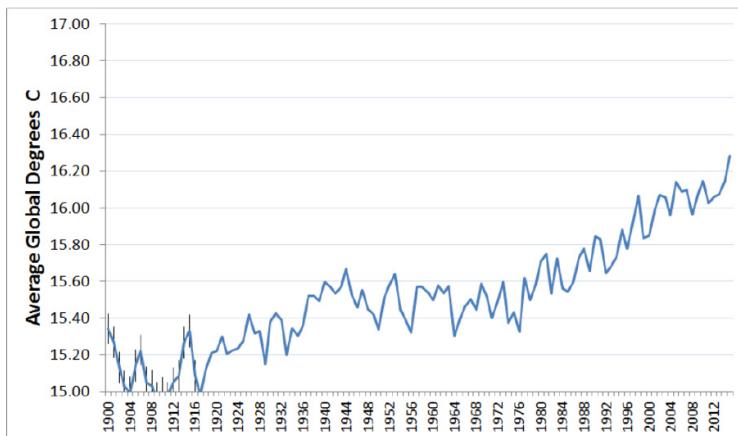
Size = radius



Size = area



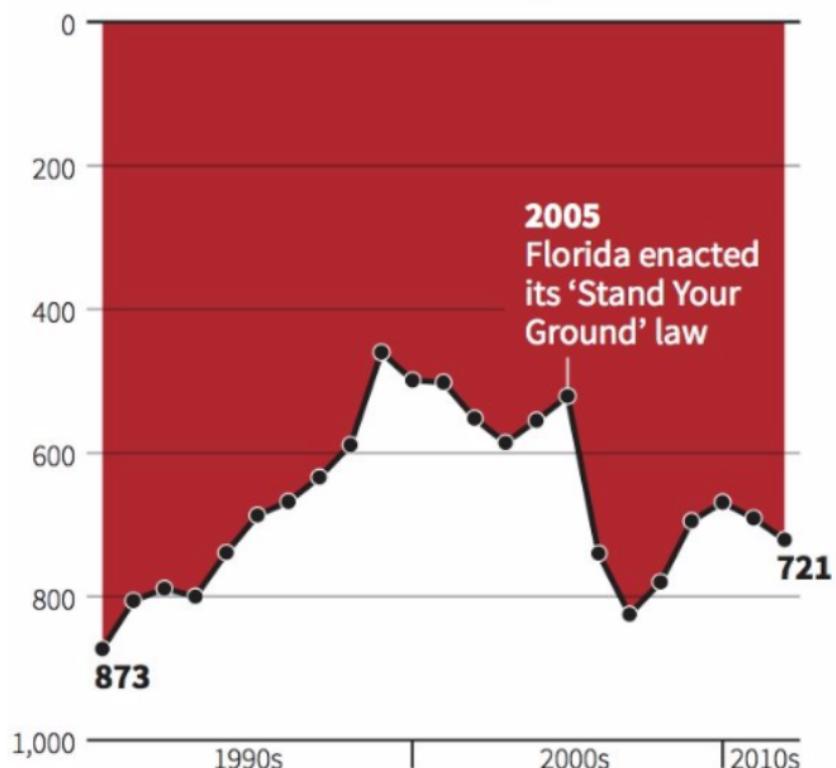




“In general, in a time-series, use a baseline that shows the data not the zero point” - [Edward Tufte](#)

Gun deaths in Florida

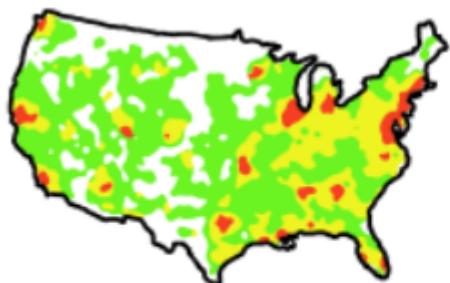
Number of murders committed using firearms



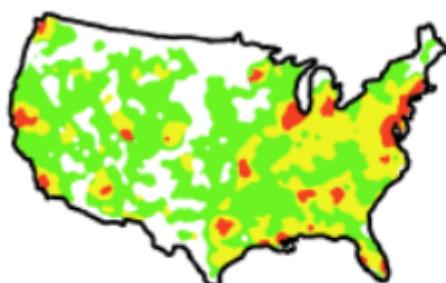
Source: Florida Department of Law Enforcement

C. Chan 16/02/2014

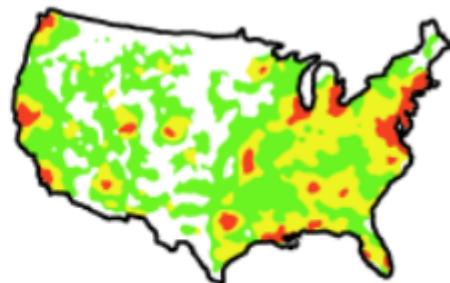
REUTERS



OUR SITE'S USERS

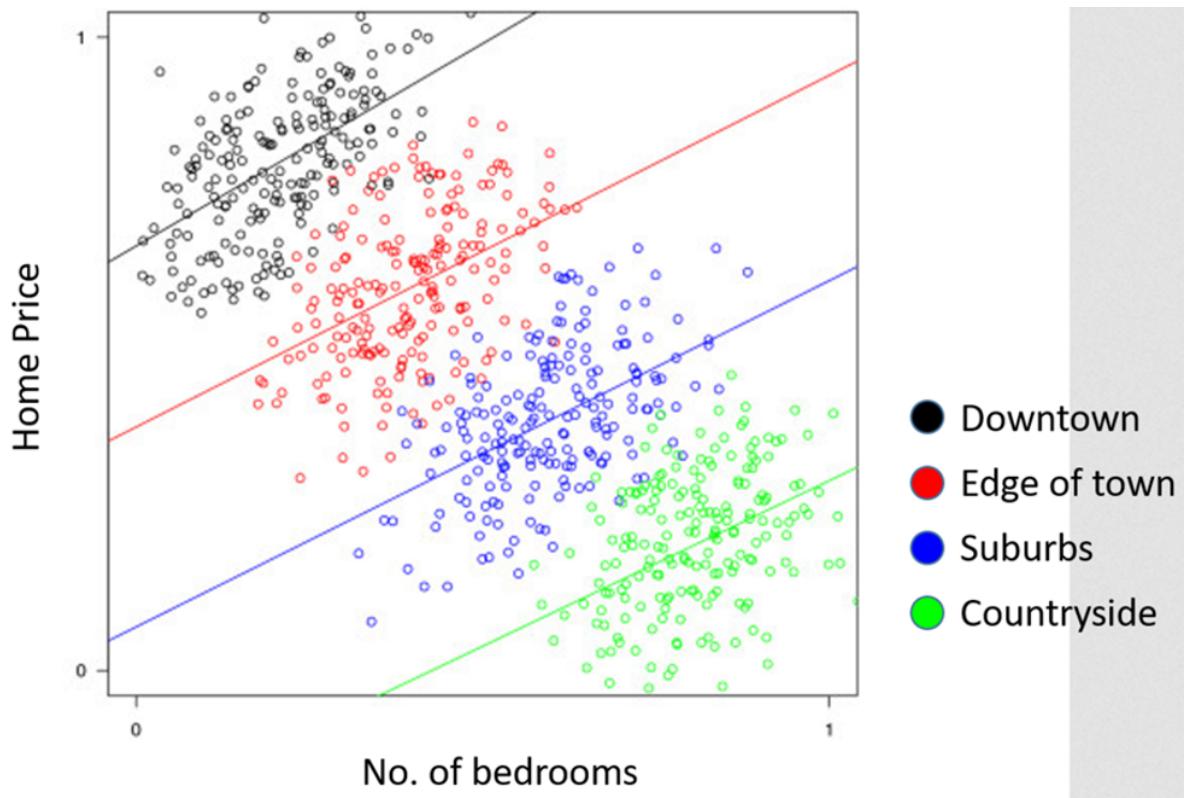


OUR SITE'S USERS



SUBSCRIBERS TO
MARTHA STEWART LIVING

Paradojas matemáticas: Simpson



Fairness

Equidad en los modelos de Machine Learning

Current Charges

<input type="checkbox"/> Homicide	<input checked="" type="checkbox"/> Weapons	<input checked="" type="checkbox"/> Assault	<input type="checkbox"/> Arson
<input type="checkbox"/> Robbery	<input type="checkbox"/> Burglary	<input type="checkbox"/> Property/Larceny	<input type="checkbox"/> Fraud
<input type="checkbox"/> Drug Trafficking/Sales	<input type="checkbox"/> Drug Possession/Use	<input type="checkbox"/> DUI/OUIL	<input checked="" type="checkbox"/> Other
<input type="checkbox"/> Sex Offense with Force	<input type="checkbox"/> Sex Offense w/o Force		

1. Do any current offenses involve family violence?
 Yes No
2. Which offense category represents the most serious current offense?
 Misdemeanor Non-violent Felony Violent Felony
3. Was this person on probation or parole at the time of the current offense?
 Probation Parole Both Neither
4. Based on the screener's observations, Is this person a suspected or admitted gang member?
 No Yes
5. Number of pending charges or holds?
 0 1 2 3 4+
 0 1 2+
6. Is the current top charge felony property or fraud?
 No Yes

Criminal History

Exclude the current case for these questions.

7. How many times has this person been arrested before as an adult or juvenile (criminal arrests only)?
5
8. How many prior juvenile felony offense arrests?
 0 1 2 3 4 5+
9. How many prior juvenile violent felony offense arrests?
 0 1 2+
10. How many prior commitments to a juvenile institution?
 0 1 2+

All Defendants			Black Defendants			White Defendants		
	Low	High		Low	High		Low	High
Survived	2681	1282	Survived	990	805	Survived	1139	349
Recidivated	1216	2035	Recidivated	532	1369	Recidivated	461	505
FP rate: 32.35			FP rate: 44.85			FP rate: 23.45		
FN rate: 37.40			FN rate: 27.99			FN rate: 47.72		
PPV: 0.61			PPV: 0.63			PPV: 0.59		
NPV: 0.69			NPV: 0.65			NPV: 0.71		
LR+: 1.94			LR+: 1.61			LR+: 2.23		
LR-: 0.55			LR-: 0.51			LR-: 0.62		

Risk scores: black vs white defendants

Number of defendants in each risk-score group



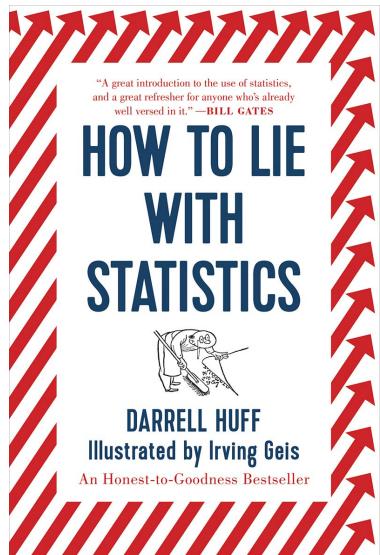
	All Defendants		Black Defendants		White Defendants	
	Low	High	Low	High	Low	High
Survived	2681	1282	990	805	1139	349
Recidivated	1216	2035	532	1369	461	505
FP rate: 32.35			FP rate: 44.85		FP rate: 23.45	
FN rate: 37.40			FN rate: 27.99		FN rate: 47.72	
PPV: 0.61			PPV: 0.63		PPV: 0.59	
NPV: 0.69			NPV: 0.65		NPV: 0.71	
LR+: 1.94			LR+: 1.61		LR+: 2.23	
LR-: 0.55			LR-: 0.51		LR-: 0.62	

¿Es el modelo justo?

En el bootcamp de Data Science aprenderás programación para explorar la data, procesarla, aplicar modelos matemáticos, deep learning, nube... Lo más importante: aprenderás a pensar

en abstracto, a plantear soluciones que un ordenador puede ejecutar

Referencias



<https://mathwithbaddrawings.com/>

¡GRACIAS! :)