# Taller de fundamentos de visualización de datos

Hola, este es el material que usé para el taller en la Tarugo#4 (octubre '19) y aunque algunas cosas se entienden bien, este contenido no está pensado para ser auto explicativo.

Cualquier pregunta a [javi@tinybird.co](mailto:javi@tinybird.co) espero que os sirva de ayuda y por lo menos os obligue a pensar un poco la próxima vez que representéis datos en pantalla.

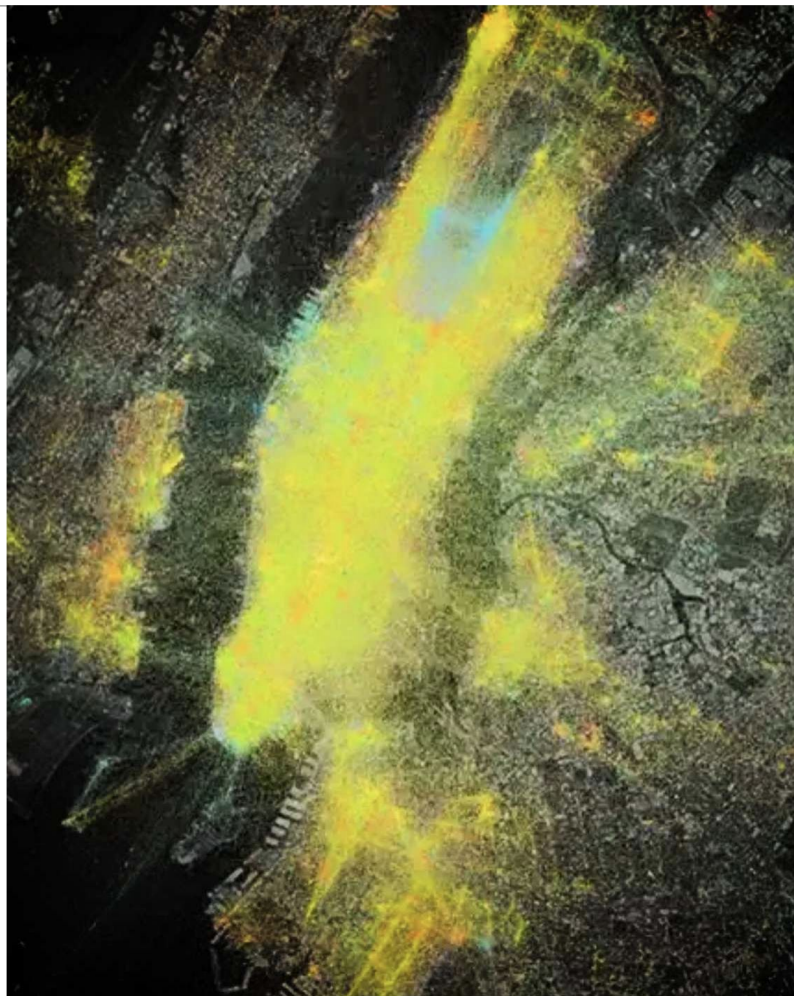Un saludo
Javi Santana
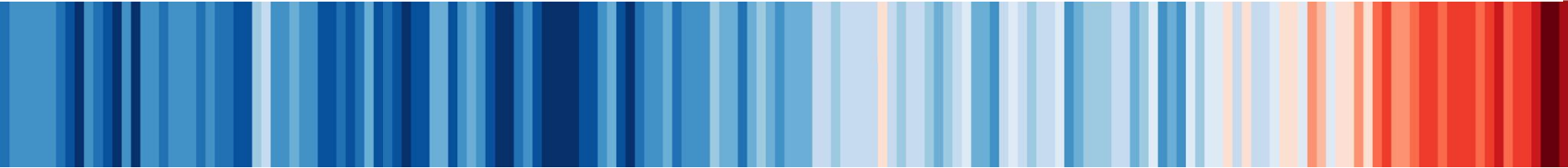
# Data Visualization

by Tinybird

# Ten Years On, Foursquare Is Now Checking In to You Even the company is still trying to figure out whether that's "cool or creepy."

*By James D. Walsh*

# The Economist

Iran's dangerous game

Lessons from a Wall Street titan

Why rent controls are wrong-headed

Goddess of the Taiwan Strait

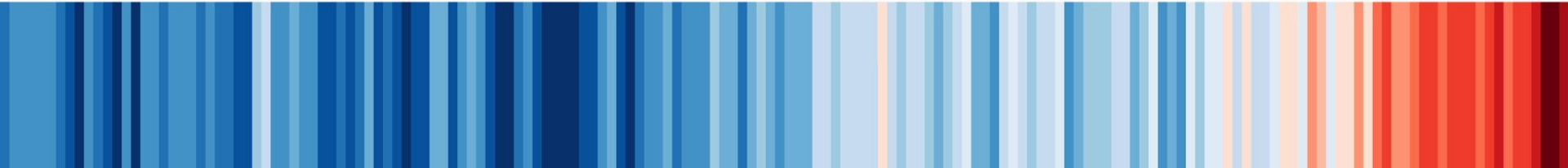SEPTEMBER 21ST–27TH 2019

# The climate issue

1850     1900     1950     2000

"I was looking for a way to communicate to audiences that aren't used to seeing graphs, or axes, or labels — things that we see day-to-day, but are complicated to them. It may look too mathematical to them, so it turns them off straight away."
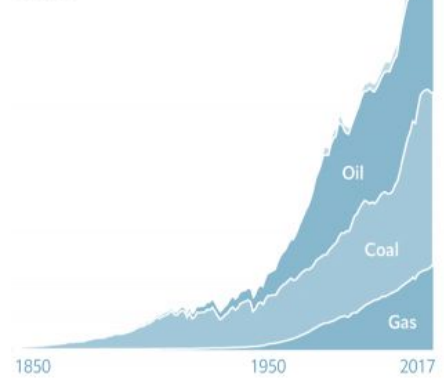— @ed_hawkins

THE DATAVIZ COMMUNITY

EFFECTIVE VISUAL COMMUNICATION

LEGENDS AND PRECISION AND THE NON-ZERO Y-AXIS DEBATE

imgflip.com

**CO₂ emissions, gigatonnes**

By fuel

Oil

Coal

Gas

1850    1950    2017

By country/region

30

China

25

United
States

20

Asia
Pacific

15

Middle
East

10

India

Americas

Africa

5

Europe

Sources: Le Quéré et al. (2018); Global
Carbon Project (GCP); Carbon Dioxide
Information Analysis Centre (CDIAC)

0

1850    1875    1900    1925    1950    1975    2000    2017

The Economist

01.

What's Data visualization?

## What's Data visualization?

## Wikipedia

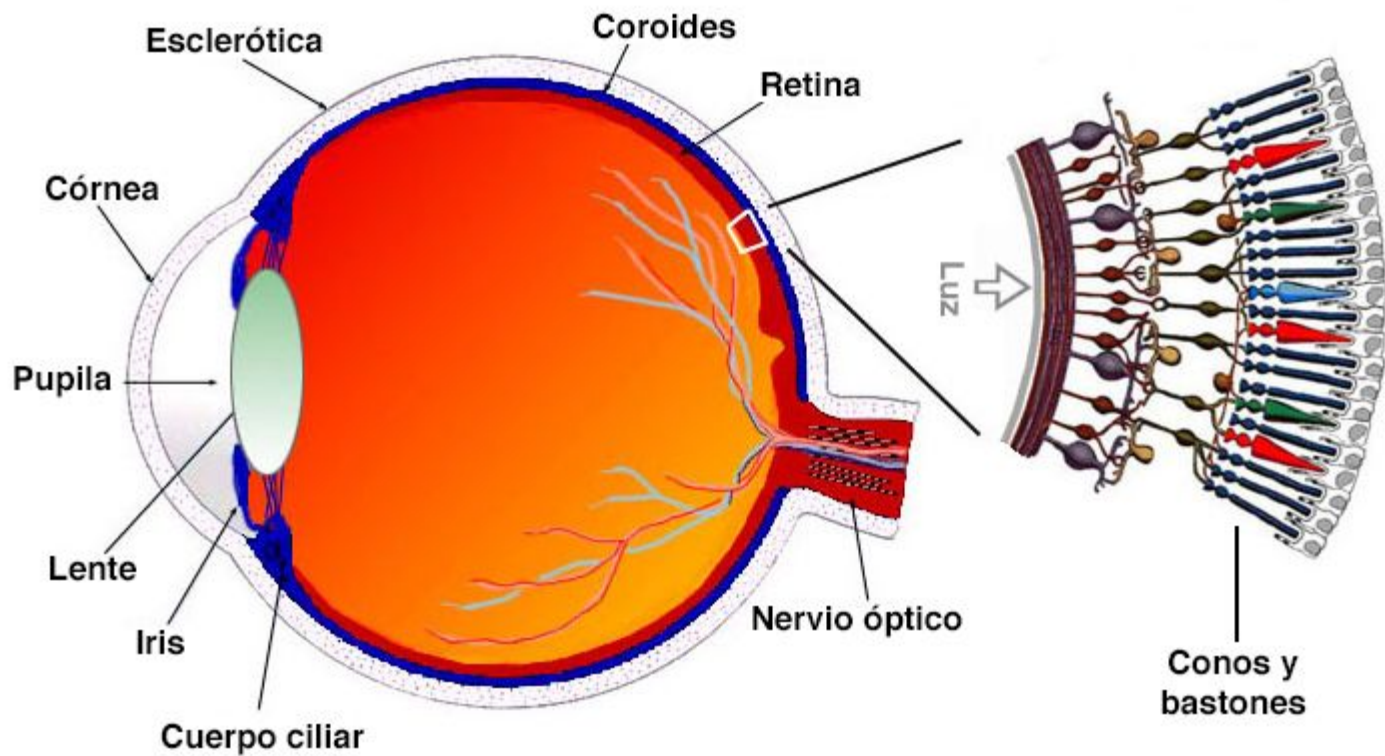Data visualization is the graphical representation of data

## What's Data visualization?

# A number is data visualization

What's Data visualization?

A number is data visualization

Esclerótica

Coroides

Córnea

Retina

Pupila

Lente

Iris

Cuerpo ciliar

Nervio óptico

Luz

Conos y bastones

Perception

**120M** bastones vs **7M** conos

# Perception II

Bastones ————————————————————————

Conos ——

# Perception II

Bastones

Conos

## Perception III

Bastones

Conos

Perception (and IV)

**Bastones**    Madrid -> Valencia
**Conos**      Madrid -> Leganés

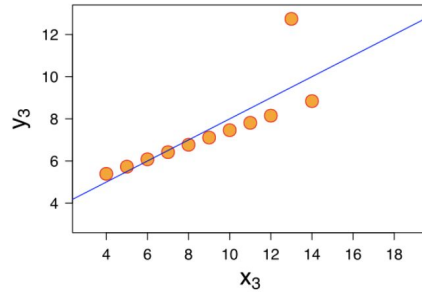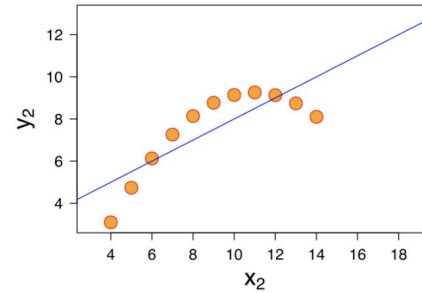# Why do we need a graphical representation if we have numbers and feelings

02.

The basis

How to represent different kind of data

Magnitudes: age, speed...
Categories: political party, football team...

**Channels:** Expressiveness Types and Effectiveness Ranks

⊕ **Magnitude Channels: Ordered Attributes**

Position on common scale

Position on unaligned scale

Length (1D size)

Tilt/angle

Area (2D size)

Depth (3D position)

Color luminance

Color saturation

Curvature

Volume (3D size)

Most

Effectiveness

Least

Same

Same

⊕ **Identity Channels: Categorical Attributes**

Spatial region

Color hue

Motion

Shape

http://hint.fm/wind/

How to represent different kind of data

- Position and size: easy to understand
- Color: not really

NOT IDEAL

HUE-BASED GRADIENT

BETTER

LIGHTNESS-BASED GRADIENT

# Using color to represent **magnitudes**

Single-Hue Sequential Color Scale

Multi-Hue Sequential Color Scale

ONE HUE

TWO HUES

# Using color to represent **categories**



The best of both worlds.

Partitions

- Differentiate colors is hard so you need to make differences more explicit

# Partitions

Basics

Partitions

Statistics + Visual representation

Color is art

No matter what stupid combination you choose there is always someone who made art with those colors

Live proof

What's Data visualization?

Statistics + Visual representation

What do I want to achieve with the visualization ?

- See global patterns (big data)

- Explore the data (big data, BI)

- Measure / make decisions (dashboards)

- Show results (scientist paper)

- Understand the data (data scientist)

- Tell a story (newspapers)

- Wow factor (newspapers, marketing, usually pretty shitty)


- "Business insights" -> BULLSHIT 99%, excel graphs 99%

03

Examples

## Examples

# Bar chart



Rough cost of a 30
second commerical
during the Super Bowl
▼

**$3.75 million**

$3m

Average profit earned
in 3.5 hours (about
the length of a Super
Bowl broadcast)*

2

1

**Anheuser-
Busch**  **Pepsi**  **General
Motors**  **Walt
Disney**  **Coca-
Cola**  **Viacom**  **Comcast**  **Time
Warner**  **Hyundai**

$50 million

Money spent on Super
Bowl commercials from
2002 — 2011

100

150

200

Note: this number is determined by taking the latest fiscal year earnings for each company, dividing by the
number of hours in a year, then multiplying by 3.5.

Ritchie King | Quartz                     Data: Compiled by Factset, Kantar Media, New York Times

# Line chart



**HIV And Wealth**

HIV prevalence (percent of the population living with HIV, ages 15 - 49). Average gross domestic product (GDP) per capita from 1990 to 2009 in U.S. dollars.

Source: UNAIDS

Credit: Adam Cole, Kevin Urmacher / NPR

# Line + bar chart

No dejes que la realidad te estropee una buena visualización

# LOS FUNDAMENTOS DE LAS PENSIONES

## > Gasto en pensiones
En miles de millones de euros.

*Previsión.

| Año | Valor |
|-----|-------|
| 2005 | 79,2 |
| 2006 | 84,6 |
| 2007 | 91,4 |
| 2008 | 98 |
| 2009 | 106 |
| 2010 | 108,2 |
| 2011 | 112,2 |
| 2012 | 115,8 |
| 2013 | 121,5 |
| 2014 | 127,4 |
| 2015 | 131,6 |
| 2016 | 135,4 |
| 2017 | 139,6 |
| 2018 | 144,5 |
| 2019* | 153,8 |

## > Beneficiarios en la Seguridad Social
En número.

*Septiembre.

| Año | Valor |
|-----|-------|
| 2014 | 8.428.617 |
| 2015 | 8.461.153 |
| 2016 | 8.609.085 |
| 2017 | 8.610.495 |
| 2018 | 8.755.362 |
| 2019* | 8.862.296 |

## > Fondo de Reserva de la Seguridad Social
En millones de euros, al cierre de cada ejercicio.

| Año | Valor |
|-----|-------|
| 2014 | 41.634 |
| 2015 | 32.481 |
| 2016 | 15.020 |
| 2017 | 8.095 |
| 2018 | 5.060 |
| 2019* | 2.000 |

* Diciembre.

Fuente: Seguridad Social

Expansión

Examples

Area chart



Exports and Imports to and from DENMARK & NORWAY from 1700 to 1780.

BALANCE in FAVOUR of ENGLAND.

Line of Imports

BALANCE AGAINST

Line of Exports

Exports

Imports

The Bottom line is divided into Years, the Right hand line into L10,000 each.

Published as the Act directs, 14 May 1786, by Wm Playfair

Neele sculpt 352, Strand, London.

# Area chart + line



**Figure 5: Evolution of rents and prices for *High Airbnb Area* vs. the rest**

(a) ln(Rents)

(b) ln(Prices)- ITP

(c) ln(Prices)- Posted

(d) Airbnb Count

*Notes:* Graph plots raw averages and the appropriate confidence intervals. *High Airbnb Area* are those in the top decile of the Airbnb listings distribution in 2016.

https://nadaesgratis.es/admin/el-impacto-de-airbnb-en-el-mercado-de-vivienda-de-barcelona/attachment/191023-airbnb

# Histogram

**Player heights**

# Heatmap (be careful with heatmaps)

**Since the Industrial Revolution, hundreds of billions of tonnes of carbon have entered our atmosphere.**

Atmospheric carbon compared with a 20th-century average, parts per million

Less carbon ▭▬ More carbon

2018
405 ppm

1844
1843
1842
1841

1860   1880   1900   1920   1940   1960   1980   2000

1840
285 ppm

**The surface of the sea is getting hotter.**

Sea-surface temperature compared with a 20th-century average, degrees Fahrenheit

Cooler ▭▬ Hotter

2018
65.2

1854
64

1860   1880   1900   1920   1940   1960   1980   2000

## Examples

### Boxplot



**OUTLIER** More than 3/2 times of upper quartile

**MAXIMUM** Greatest value, excluding outliers

**UPPER QUARTILE** 25% of data greater than this value

**MEDIAN** 50% of data is greater than this value; middle of dataset

**LOWER QUARTILE** 25% of data less than this value

**MINIMUM** Least value, excluding outliers

**OUTLIER** Less than 3/2 times of lower quartile

# Boxplot



Age distribution of Olympic Athletes by Sport and Gender: All-time
Female = Pink, Male = Blue, Both = Green

Examples

# Scatterplot

Power generation by source (2000-2018)

■ coal　■ other fossil fuels　□ nuclear
■ renewables

2000　2018

Austria　Belgium　Bulgaria　Croatia
Cyprus　Czech Rep.　Denmark　Estonia
Finland　France　Germany　Greece
Hungary　Ireland　Italy　Latvia
Lithuania　Luxembourg　Malta　Netherlands
Poland　Portugal　Romania　Slovakia
Slovenia　Spain　Sweden　UK

Source: Sandbag Climate Campaign

https://twitter.com/ghensel/status/1133350342819815426

Maps

PUMP

CARSLI

BATEMANS B

GREAT MARLBOROUGH STREET

PORTLAND STREET

WARDOUR MEWS

FRITH

ARGYLL PLACE

WORK HOUSE

PORTLAND MEWS

W I C K

St. ANNS LANE St. ANNS COURT

RICHMOND BUILDINGS

GREEN DRAGON YARD

PUMP

MARSHALL ST

BENTINCK STREET

BREWERY YARD

RICHMOND MEWS

QUEEN

LITTLE MARLBOROUGH ST

TYLER COURT

COWNES CT

EDWARD STREET

DUCK LANE

CARNABY

PETS PLACE

DUFOURS PLACE

BROAD

STREET

HAM SQUARE

MEARDS STREET

TYLER STREET

TOUBERTS PL

KING

CROSS ST

SOUTH ROW

BROAD

PUMP

CANNON ROW ST

BREWERY

MEARDS COURT

MEARDS COURT

CROSS STREET

NEW STREET

COCK CT

MAIDENHEAD COURT

MACKENS COURT

MALLOW YARD

SILVER

STREET

PETER STREET

LITTLE

DEAN

STREET

OLD COMPTON

CHAPEL PLACE

BRIDLE STREET

WINDMILL

HUSBAND ST

HOPKINS

PULTENEY COURT

GREEN COURT

WALKERS COURT

RICHMOND ST

LD BURLINGTON MEWS

BEAK ST

UP JOHN ST

UP JAMES ST

GREAT PULTENEY STREET

WILLIAM AND MART YARD

LITTLE PULTENEY

PUMP

NEW BURLINGTON STREET

ST JAMES ST

GOLDEN SQUARE

QUEENS HEAD COURT

GT CROWN CT

GREAT CROWN CT

NT CROWN COURT

KING STREET

PUMP

UP RUPERT ST

BLOGH YARD

GEORGE

# The Facebook Offering: How It Compares

Find a company [            ]

**Company value**
In billions of today's dollars

100 —

80 —

60 —

40 —

20 —

10 —

1 —

0.1 —

**Facebook**

This is the same chart on a logarithmic scale. With this scale, percentage increases and decreases are comparable.

1980          1985          1990          1995          2000          2005          2010

**Year of I.P.O.**

Angel Island, police search location

Paddle found here

Probable escape location

Escape from Alcatraz

Alcatraz

Golden Gate Bridge

Time of escape
- 20:00
- 21:00
- 22:00
- 23:00
- 00:00
- 01:00
- 02:00
- 03:00
- 04:00

## More animated maps

https://flowingdata.com/2015/12/15/a-day-in-the-life-of-americans/
http://javisantana.com/kotar/
https://team.carto.com/u/javi/me
https://javi.carto.com/builder/870c41d2-fc35-11e3-83a9-0edbca4b5057/embed

En realidad puedes hacer lo que quieras
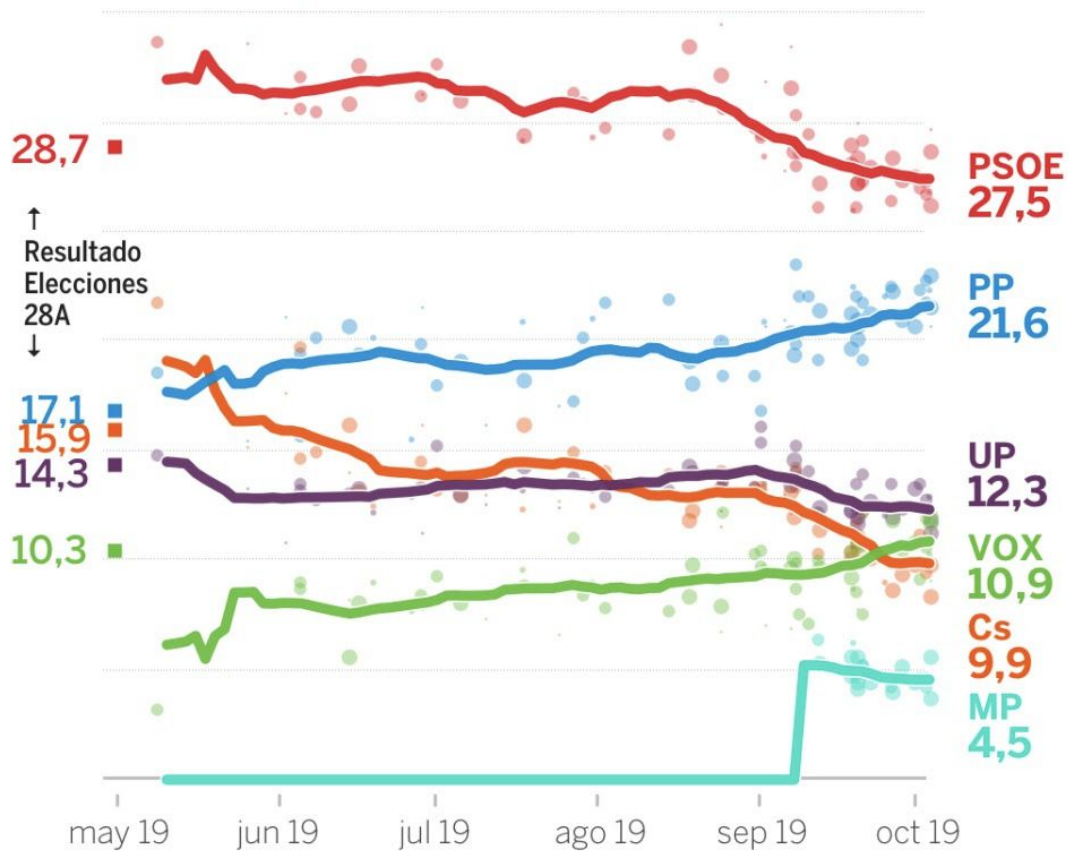
# Estimación de voto según cada encuesta, y **promedio** a partir de todas ellas



28,7 ■

↑
Resultado
Elecciones
28A
↓

17,1 ■
15,9 ■
14,3 ■

10,3 ■

**PSOE**
27,5

**PP**
21,6

**UP**
12,3

**VOX**
10,9

**Cs**
9,9

**MP**
4,5

may 19    jun 19    jul 19    ago 19    sep 19    oct 19

http://www.evolutionoftheweb.com/

# Dissecting a Trailer: The Parts of the Film That Make the Cut

How scenes from five of the nine best picture nominees were reassembled to promote the films.

## Silver Linings Playbook

"Silver Linings Playbook" follows the standard model for trailers, according to Bill Woolery, a trailer specialist in Los Angeles who once worked on trailers for movies like "The Usual Suspects" and "E.T. the Extra-Terrestrial." While introducing the movie's story and its characters, the trailer largely follows the order of the film itself.



Start of trailer →

0 sec    30 sec    60 sec    90 sec    120 sec

**Beginning of film**

**Middle**

**End**

Not in film

The trailer's opening shot — an image of the family's home — appears near the end of the film, but there are similar shots near the beginning of the movie.

A handful of very short shots are never seen in the film, although most are shown from alternate camera angles.

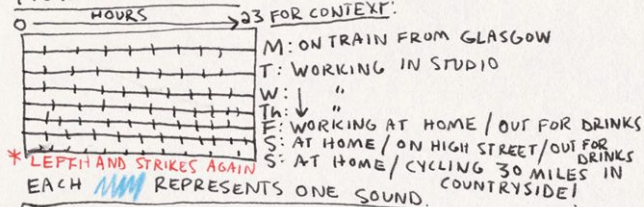Shots that accompany the main actors' names are also shown out of order.

# DEAR DATA — WEEK 32

## A WEEK OF SOUNDS

ABOUT THE DATA: EVERY HOUR I WAS AWAKE
I TRACKED THE SOUNDS I HEARD AROUND ME
(GENERALLY ON THE TOP OF THE HOUR, SOMETIMES
LATER IF I FORGOT)

HOW TO READ IT:

0 ———— HOURS ————→ 23   FOR CONTEXT:

M: ON TRAIN FROM GLASGOW
T: WORKING IN STUDIO
W: ↓ "
Th: ↓ "
F: WORKING AT HOME / OUT FOR DRINKS
S: AT HOME / ON HIGH STREET / OUT FOR DRINKS
S: AT HOME / CYCLING 30 MILES IN COUNTRYSIDE!

* LEFT HAND STRIKES AGAIN

EACH ⋙ REPRESENTS ONE SOUND.

SOUND TYPES ARE ORGANISED AS FOLLOWS:

① 'ORGANIC' SOUNDS ( SOUNDS CREATED BY PEOPLE,
ANIMALS, OR NATURE)

A SOUND MY HUSBAND MADE: SPEAKING,
DRINKING COFFEE, HUMMING, FIXING BIKE, ETC.*

PEOPLE'S MOVEMENT SOUNDS:
FOOTSTEPS, EATING, PUSHING CHAIRS BACK
( SOUNDS OF ACTIVITY )

PEOPLE'S VOICES

CLANKING BOTTLES + CUTLERY

ROLLING WHEELS: SUITCASES, CARTS, ETC.

RUSTLING, SHUFFLING NEWSPAPERS

CRASHES, BANGS, + RATTLES

BIRDSONG

WIND IN TREES

② 'MACHINE' SOUNDS:

RUNNING WATER

APPLIANCES RUNNING:
RADIATOR / BOILER, WASHING MACHINE, POWER TOOL,
COFFEE MACHINE, ETC.

RECORDED MUSIC

LAPTOP

OUR FILM / TV

MOTOR VEHICLE

NEIGHBOURS' LOUD TV + MUSIC

TRAINS

AIRPLANE

THE HUM OF ELECTRICITY

③ UNUSUAL SOUNDS: ⋙ STEEL DRUM BAND,
ARCADE GAMES, CHURCH BELLS, A HORSE!

* ALSO INCLUDES SOME RIDICULOUSLY STUPID SOUNDS HE WAS MAKING TO MESS W/ THE DATA

FROM:
S POSAVEC ▨▨▨▨
LONDON ▨▨▨▨
UK ▨▨▨

TO:
GIORGIA LUPI
▨▨▨▨▨▨▨▨
BROOKLYN, NY ▨▨▨
USA
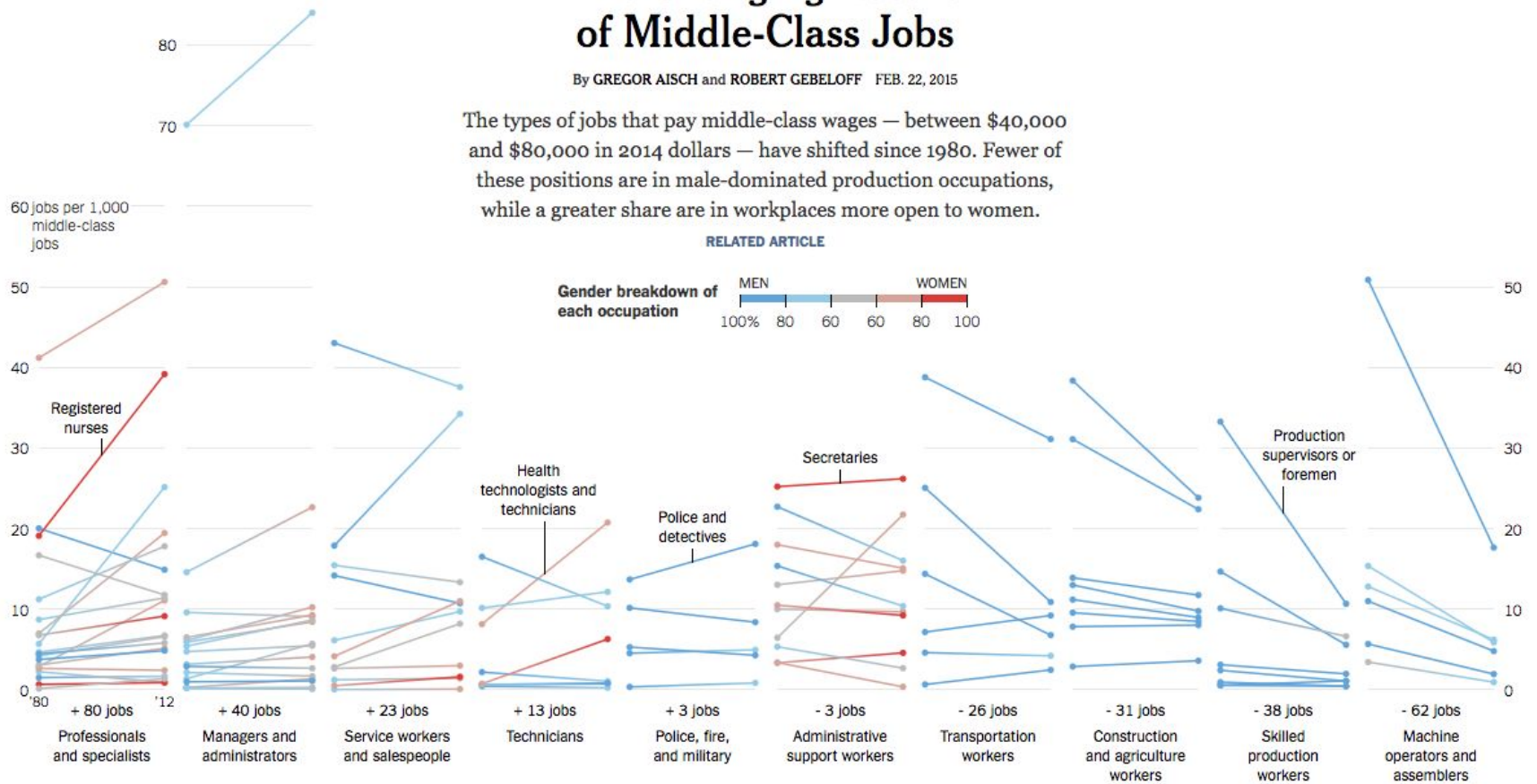
BY AIR MAIL
*par avion*
Royal Mail®

# The Changing Nature of Middle-Class Jobs
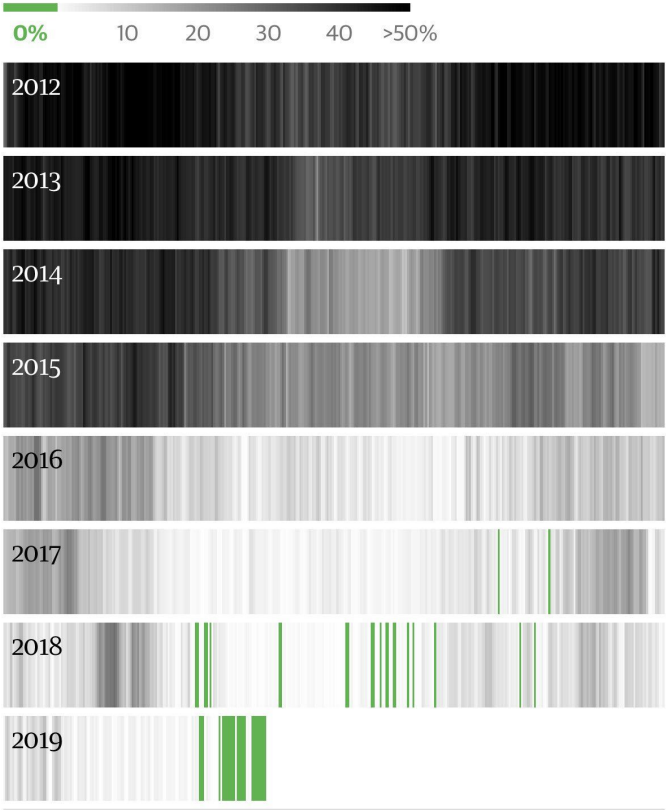
By **GREGOR AISCH** and **ROBERT GEBELOFF**   FEB. 22, 2015

The types of jobs that pay middle-class wages — between $40,000 and $80,000 in 2014 dollars — have shifted since 1980. Fewer of these positions are in male-dominated production occupations, while a greater share are in workplaces more open to women.
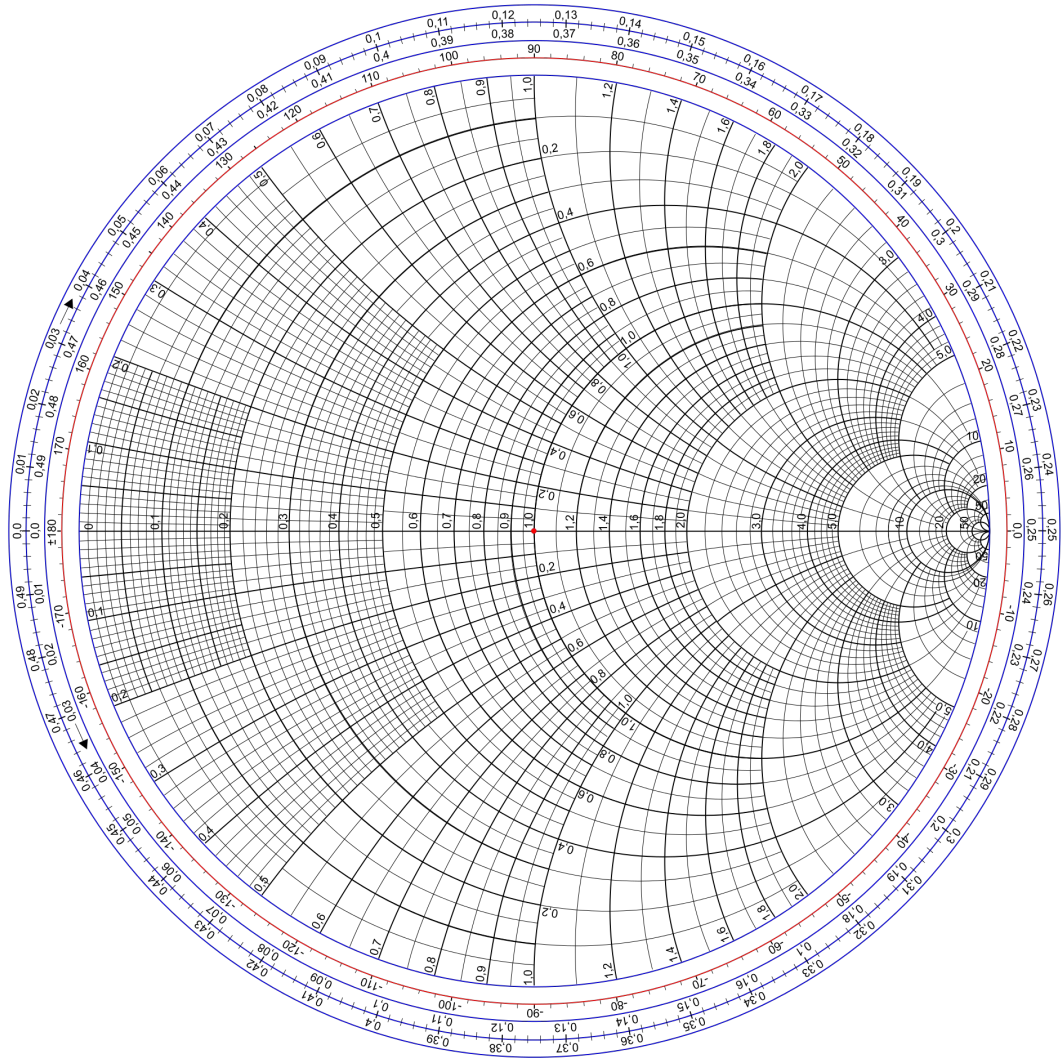
RELATED ARTICLE



Gender breakdown of each occupation — MEN 100% 80 60 | 60 80 100 WOMEN

80
70
60 jobs per 1,000 middle-class jobs
50
40
30
20
10
0

'80    '12

Registered nurses

Health technologists and technicians

Police and detectives

Secretaries

Production supervisors or foremen

| + 80 jobs | + 40 jobs | + 23 jobs | + 13 jobs | + 3 jobs | - 3 jobs | - 26 jobs | - 31 jobs | - 38 jobs | - 62 jobs |
|---|---|---|---|---|---|---|---|---|---|
| Professionals and specialists | Managers and administrators | Service workers and salespeople | Technicians | Police, fire, and military | Administrative support workers | Transportation workers | Construction and agriculture workers | Skilled production workers | Machine operators and assemblers |

**Britain is rapidly phasing out coal**

Daily share of Britain's power generated by burning coal

0%   10   20   30   40   >50%

2012

2013

2014

2015

2016

2017

2018

2019

https://twitter.com/EmmaFidler/status/1132347203031326722
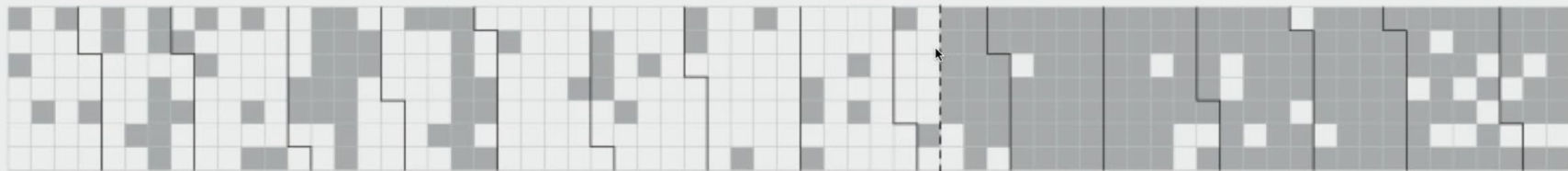
## Picture days

*Before my son was born*    *After*

2013    2014

04

Conclusions

## Tufte

- Less is usually more (ink-data ratio)
- Graphical excellence consists of complex ideas communicated with clarity, precision, and efficiency
- Keep it proportional! "Lie Factor = *size of effect shown in graphic* divided by *size of effect in data*"
- You don't have to use a graphic when there isn't much data — a table is often better.
- Pie charts are useless. Period. DO NOT USE PIE CHARTS.

Conclusions

## Santana

- La visualización de datos normalmente sirve para saber qué preguntas hacer, no para resolverlas
- "All non-trivial abstractions, to some degree, are leaky" (Joel)
- No hay una forma buena de visualizar algo
- No te dejes engañar por las pintas de una visualización
- A hacer visualización se aprende haciendo visualizaciones (y copiando y probando)
- No quieres ser este tipo de persona:

   "No dejes que la realidad te estropee una buena visualización"

05.

Must reads

Must reads

- Visual display of quantitative information (Tufte)
- @flowingdata
- @MKrzywinski
- @lisacrost
- Dataviz project https://datavizproject.com
- https://nadaesgratis.es data analysis + vis

06.

HANDS ON! (Santana's way)

## About the visualization

This visualization could be done with a library in a few minutes, we are trying to show how easy is to visualize data with no tools.

It might be too easy for experienced developers. The code is not the important part.

## NO2 in Madrid

- Lots of public data in a reusable format
- This is a big problem, we are all breathing this thing right now (72 ug/m^3 today)

# The process

- Get Data -> Visualize

# The process - the data

- Understand the source data: format, meaning and so on

- Understand the actual data: histograms, max, min...

- Clean the data

- Prepare the data to be visualized

    - Depending on the data size this step is mandatory

## Some [random] data tools

- Cleaning: Trifacta // Open refine // Python // R // …

- Prepare data: Any database SQL database (Postgres) / CSVKit / Python…

- Understand the data: excel // pandas…

- Serve: S3 // github // datasette // carto // tinybird

Santana's choice:

- Python // clickhouse // vega (tinybird)

# The process - the visualization

- Decide what we want to show

- Go for a visualization type, remember

    - See global patterns

    - Explore the data

    - Measure / make decisions

    - Tell a story

    - Wow factor

- Iterate 100 times

## Some [random] vis tools

- Vega, D3, observable

- Pandas (mathplotlib)

- Datastudio/tableu/qlik/....

Santana's choice:

- Vanilla js, d3, pandas

These charts show the percentage change of 19 technology stocks since December 2008. It's easy to spot similarities: market events that affected multiple stocks.

Si has vivido

5 años cerca de

Gran Vía

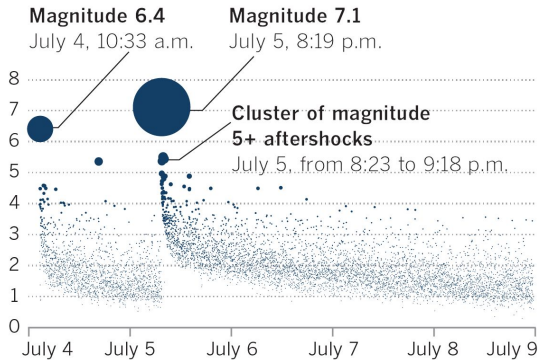has respirado el equivalente a

Comer 3.4kg de carbonilla

# What kind of vis?

# What kind of vis?

## Ridgecrest aftershocks tapering off

Thousands of aftershocks have followed two larger earthquakes near Ridgecrest last week. They've followed a predictable, decreasing pattern.
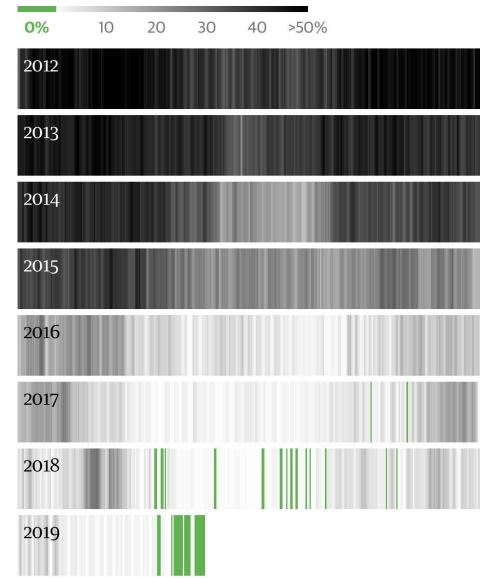
**Magnitude 6.4**
July 4, 10:33 a.m.

**Magnitude 7.1**
July 5, 8:19 p.m.

**Cluster of magnitude 5+ aftershocks**
July 5, from 8:23 to 9:18 p.m.

Points show earthquakes between 10 a.m. July 4 and 10 a.m. July 9 near Ridgecrest, Calif. The gap in aftershocks following the magnitude 7.1 is because instruments are unable to detect smaller earthquakes following a main shock.

Sources: U.S. Geological Survey, Caltech, Times Reporting

Chris Keller / Los Angeles Times

## Britain is rapidly phasing out coal
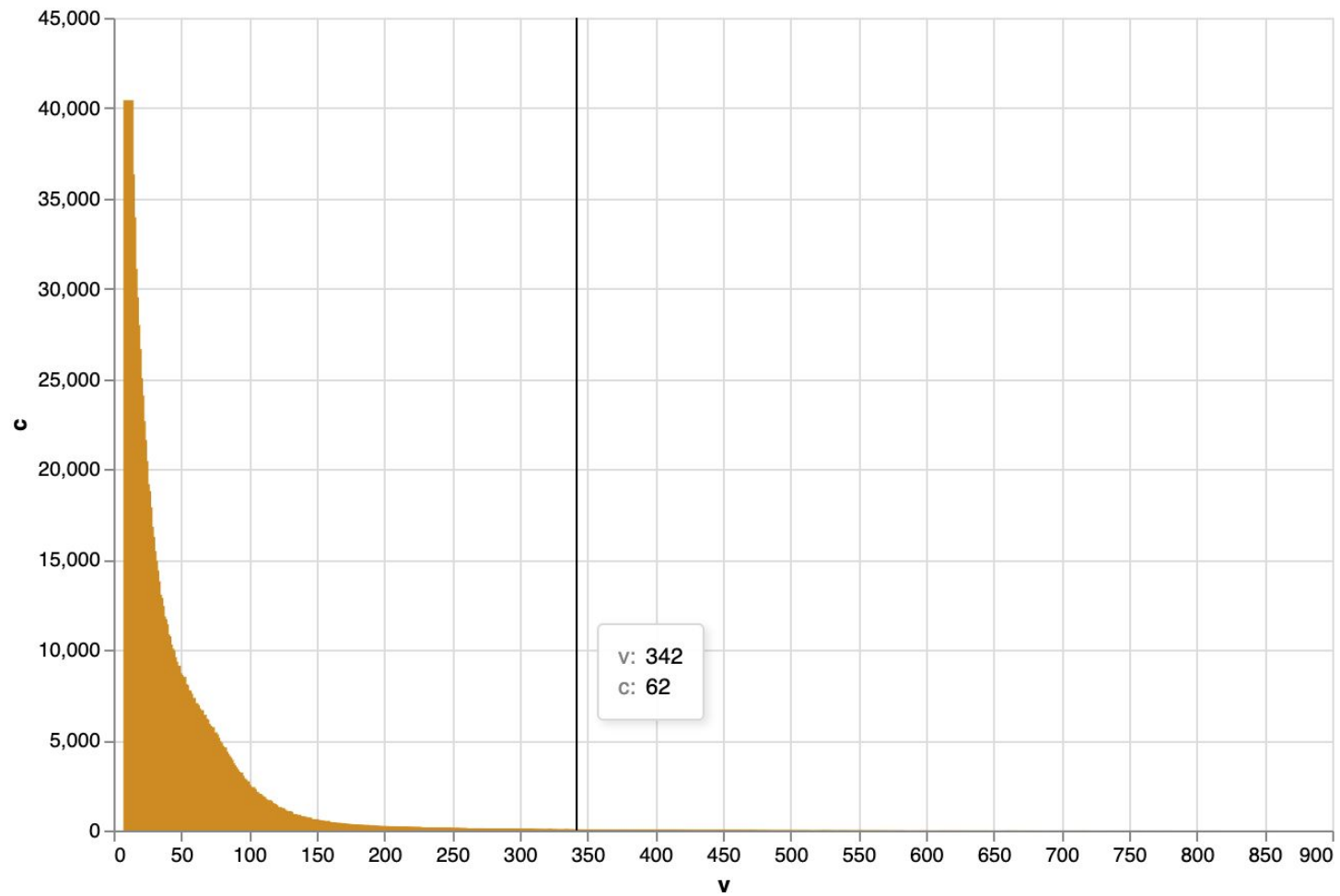
Daily share of Britain's power generated by burning coal

0%   10   20   30   40   >50%

2012

2013

2014

2015

2016

2017

2018

2019

http://javisantana.com/**data_vis_workshop**/

## Steps

- JSON

- Render a square

- Render the data

- Color

- More on colors

- Text
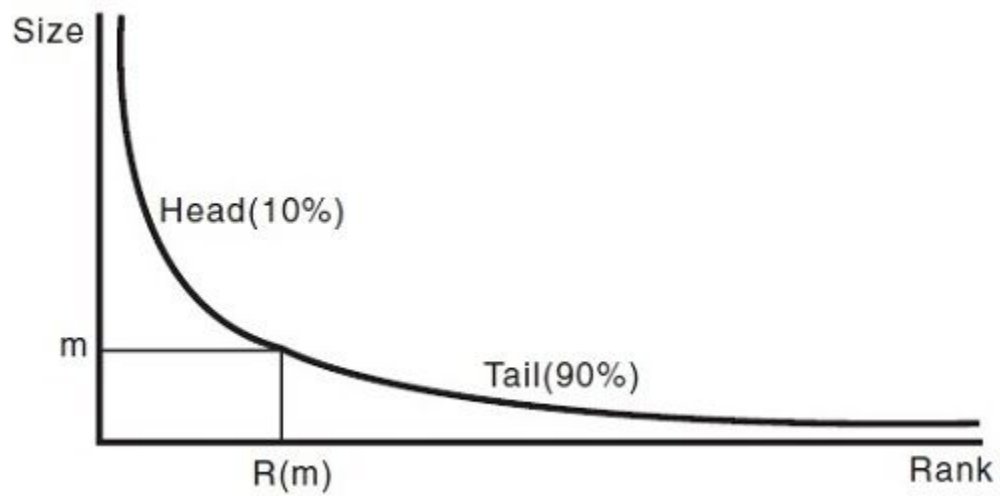
# Protocolo anticontaminación

- Alerta: 180

- Aviso: 200
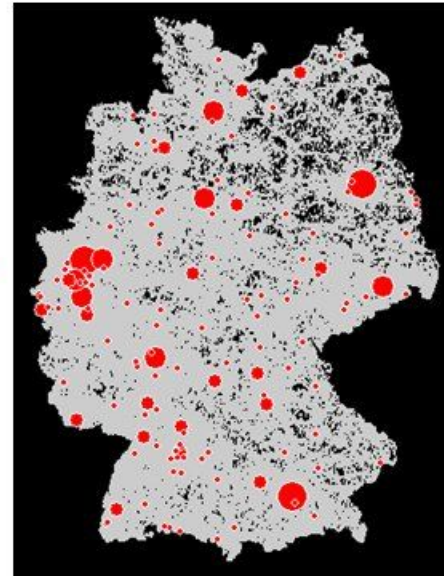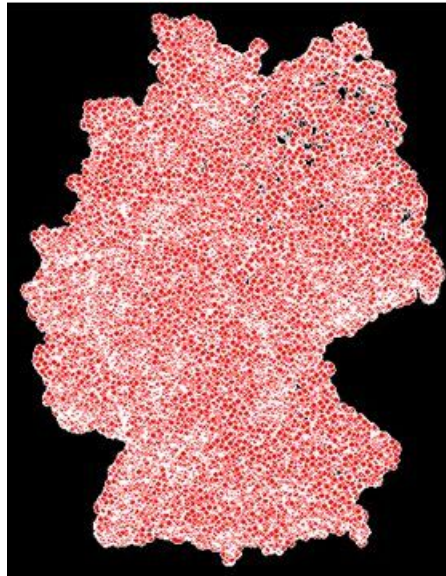
# Head/tails

Head/tails

- [34,65,91,123,169,237,321]

Thanks - QA
@javisantana
@tinybirdco

Hands on video from data to vis

In this video I don't work on the same visualization but the process is the same

https://www.youtube.com/watch?v=V1nbigUstGA