

โครงการอบรมหลักสูตร Hand-on Data Science and Machine Learning

Data Visualization

Navavit Ponganant
Data scientist
Government Big Data Institute (GBDi)

Definition

Data visualization is the graphical representation of information and data. By using visual elements like charts, graphs, and maps, data visualization tools provide an accessible way to see and understand trends, outliers, and patterns in data.

Data visualization beginner's guide: a definition, examples, and learning resources
(<https://www.tableau.com/learn/articles/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

6200 BC

1500 BC

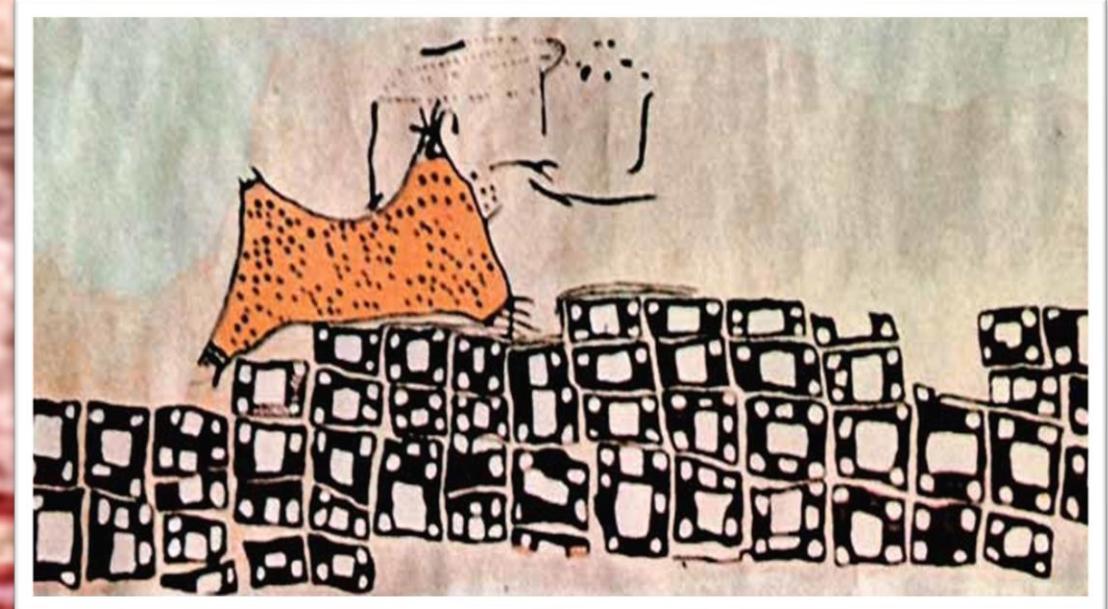
1786

1801

1855

1861

A Wall Painting at Catal Huyuk



6200 BC

1500 BC

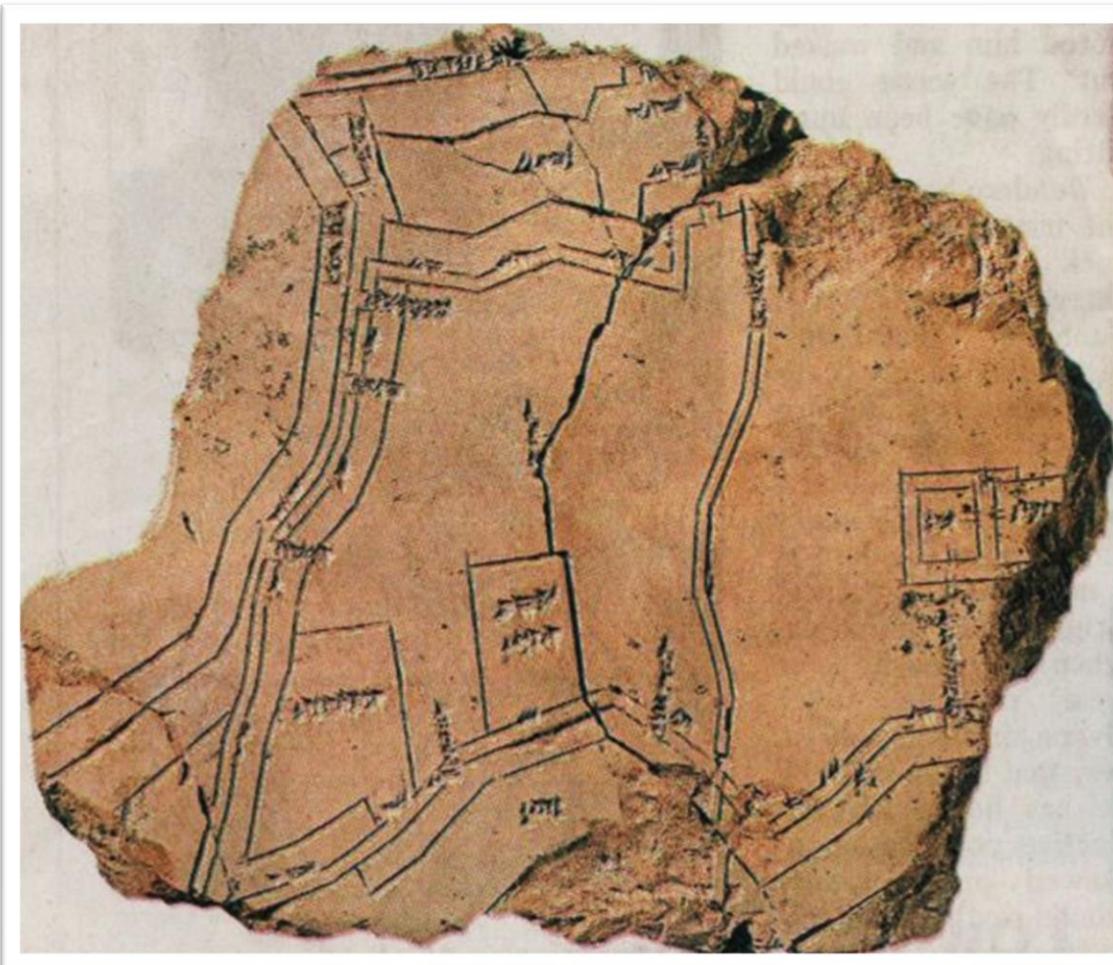
1786

1801

1855

1861

Clay tablet from Nippur, Babylonia



6200 BC

1500 BC

1786

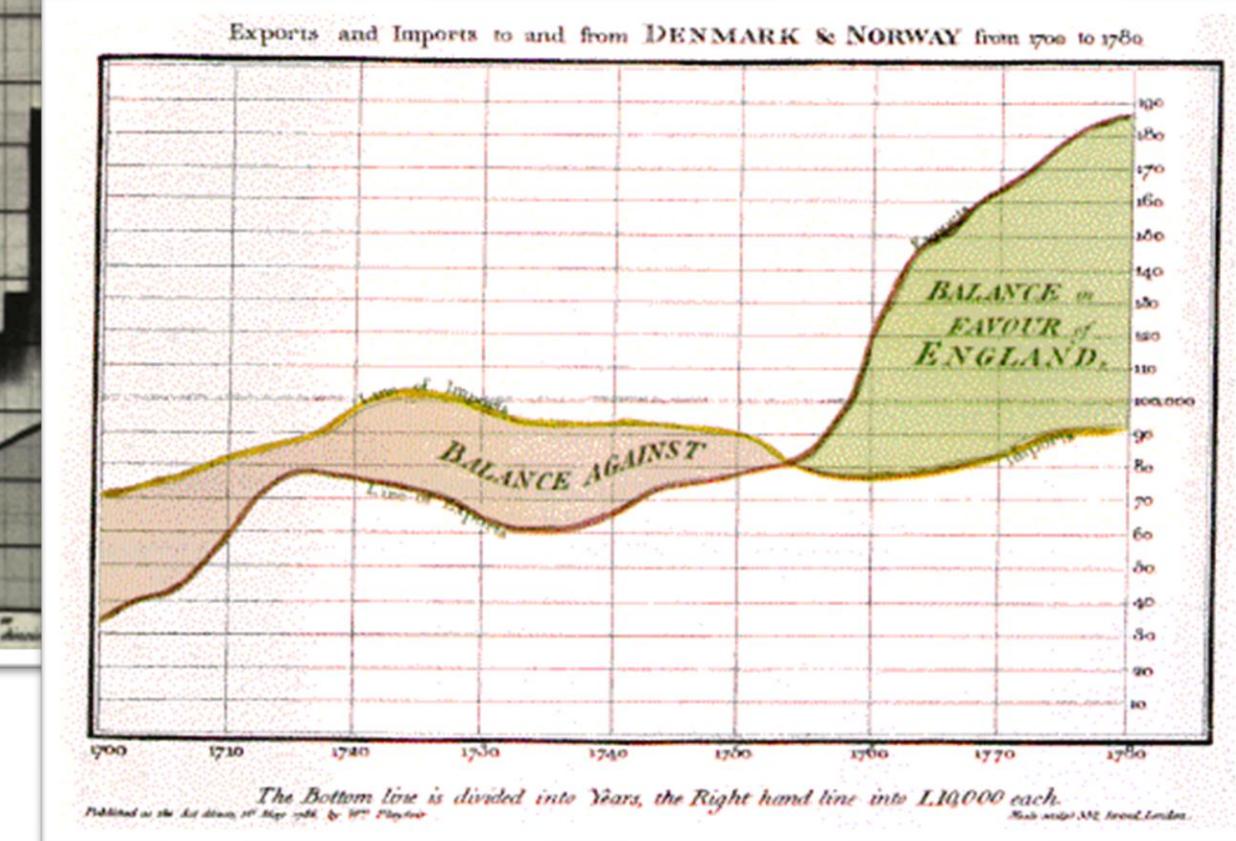
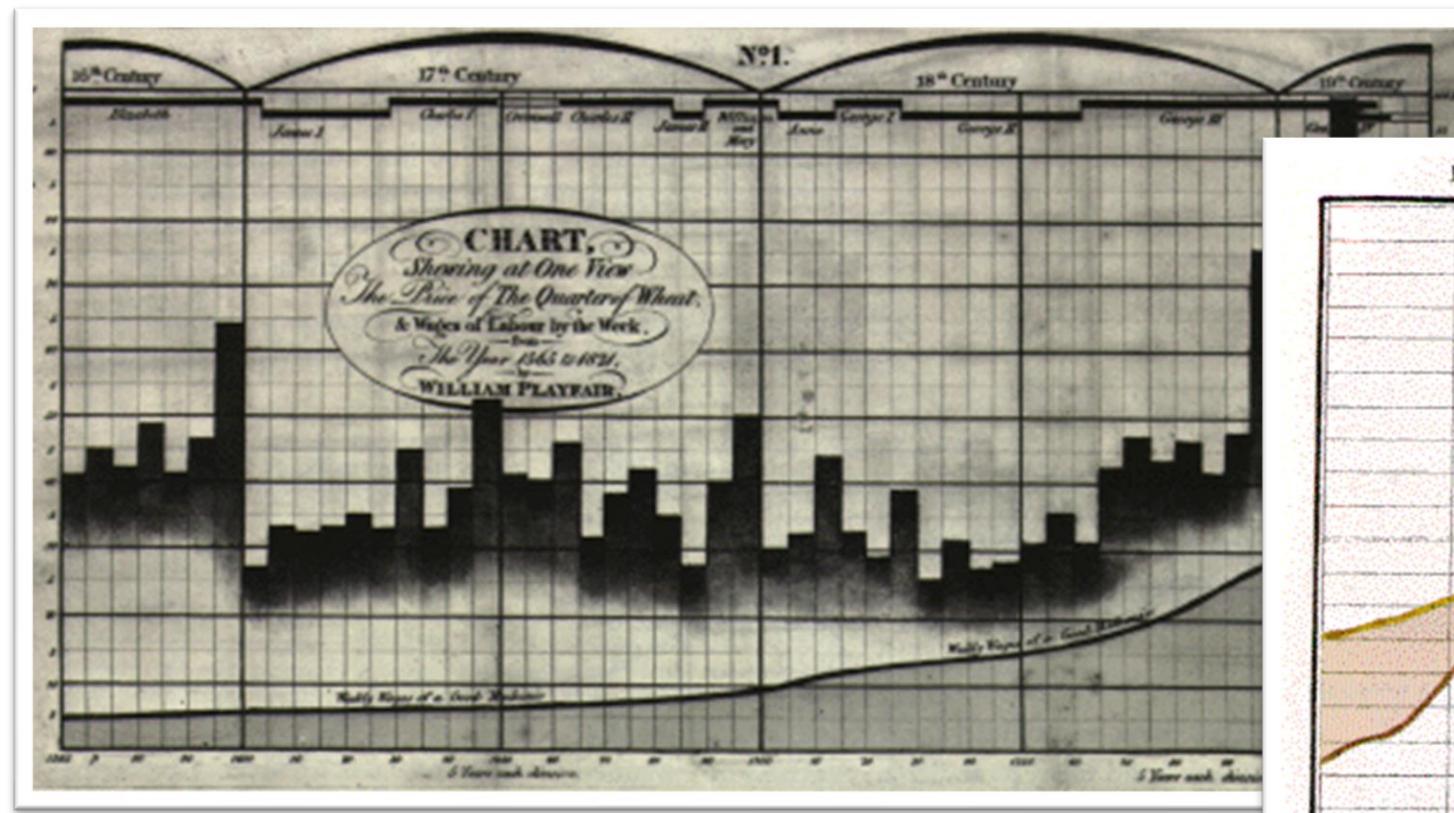
1801

1855

1861

William Playfair: The Commercial and Political Atlas (Book)

- Abstract Data Presentation



6200 BC

1500 BC

1786

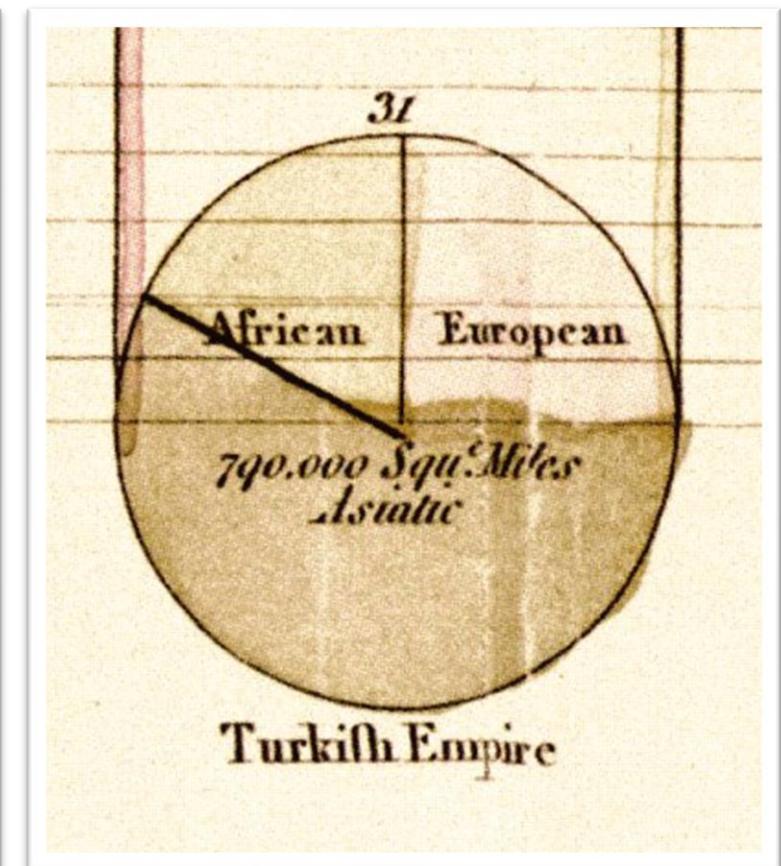
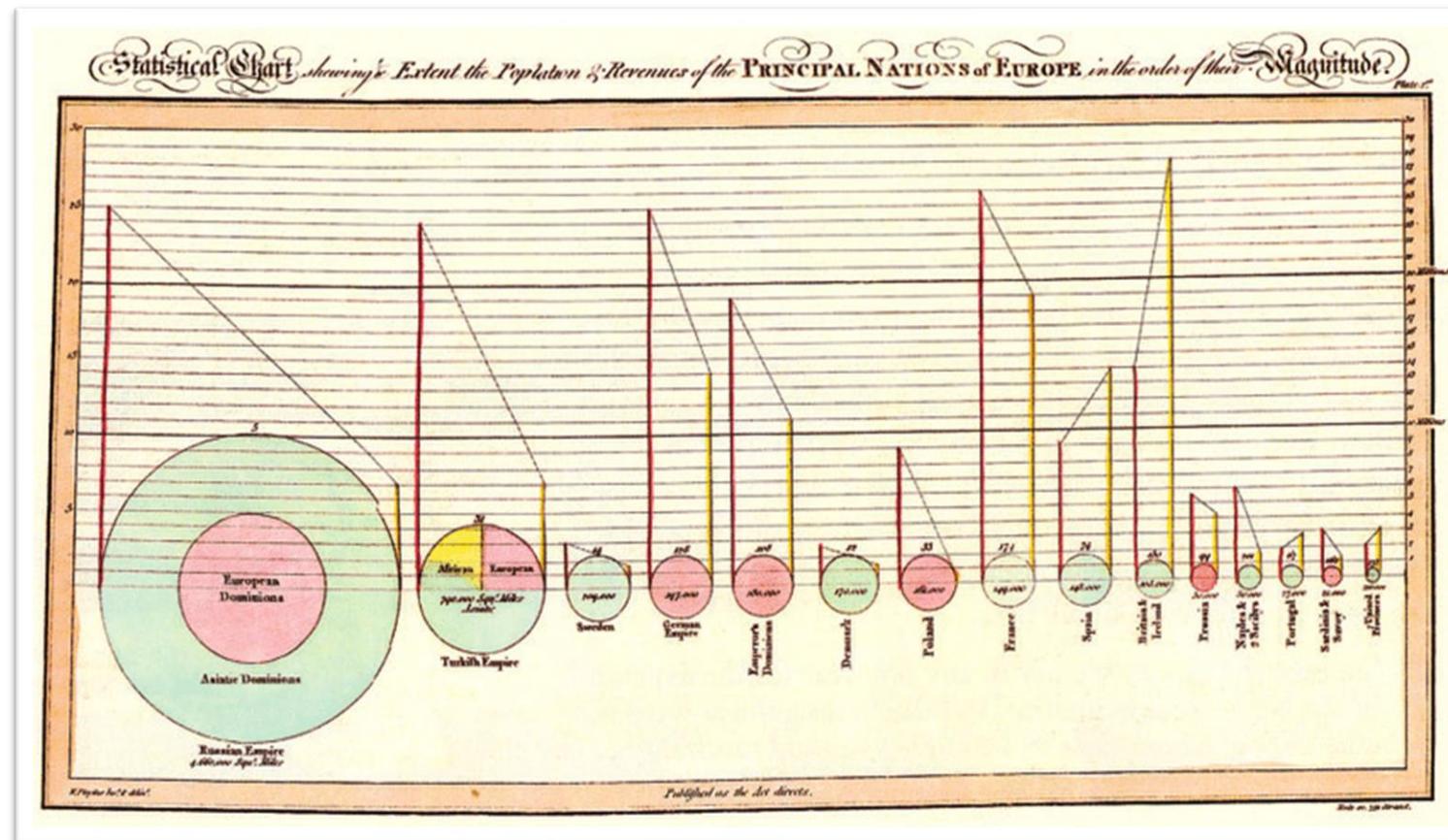
1801

1855

1861

William Playfair: A Statistical Breviary (Book)

- The earliest recorded example of a Pie Chart



6200 BC

1500 BC

1786

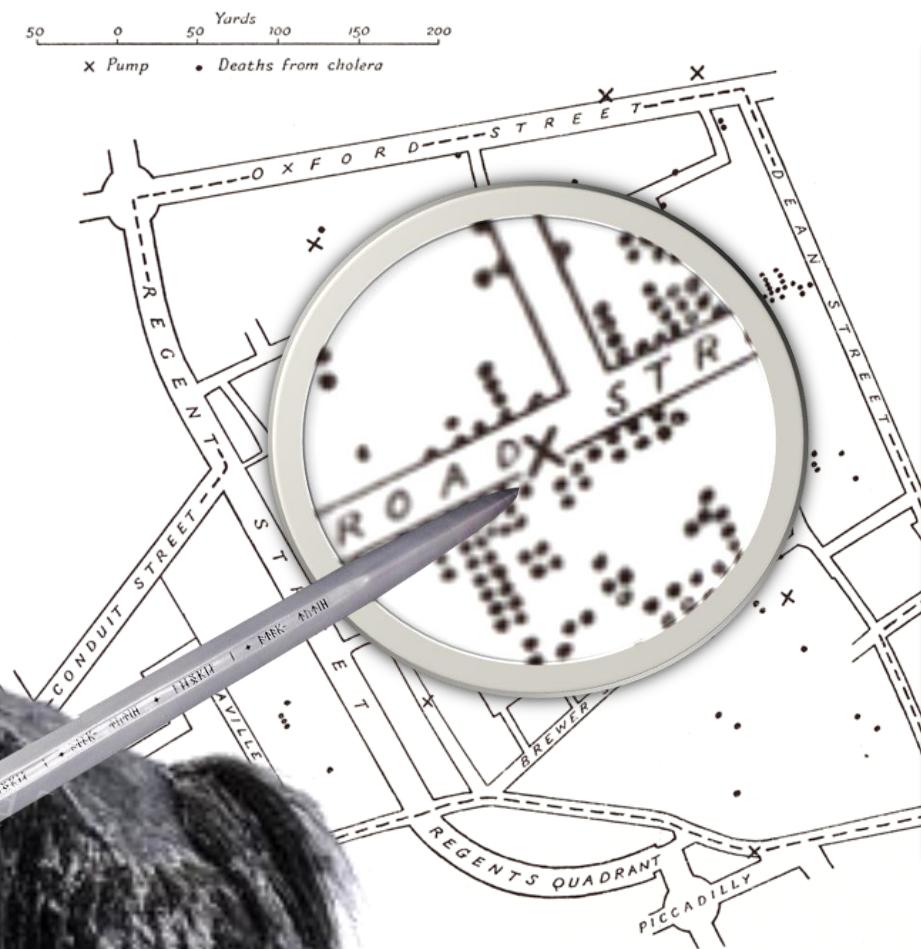
1801

1855

1861

Dr. John Snow: Statistical Map Visualization

– London Cholera Epidemic



Broad Street Pump

6200 BC

1500 BC

1786

1801

1855

1861

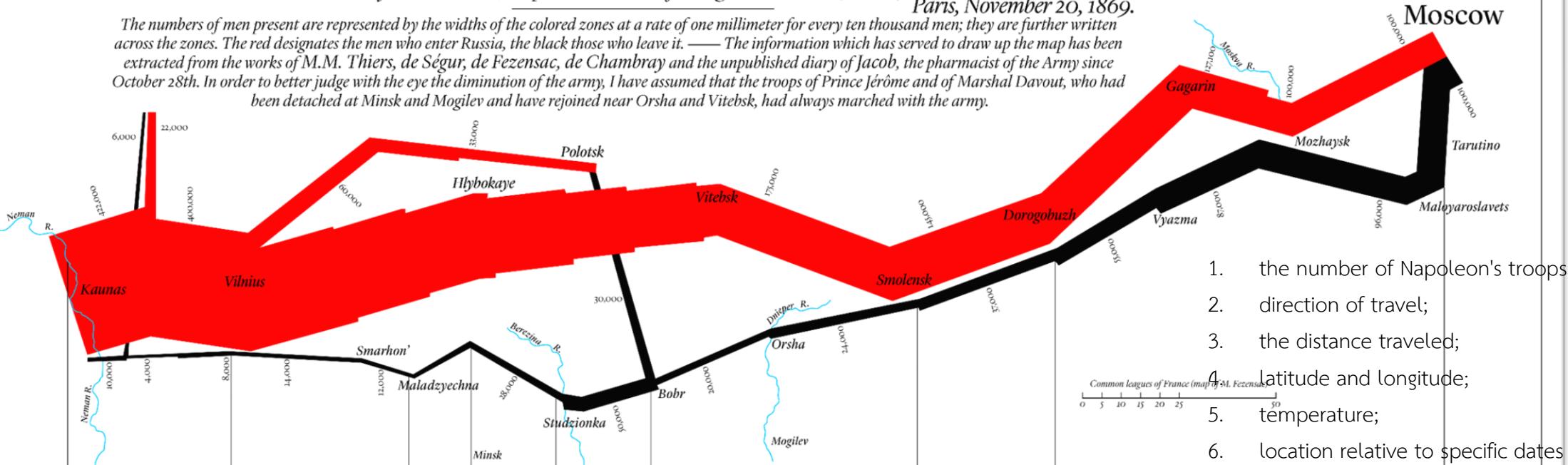
Charles Minard: Napoleon's March

Figurative Map of the successive losses in men of the French Army in the Russian campaign 1812 ~ 1813

Drawn by M. Minard, Inspector General of Bridges and Roads (retired).

Paris, November 20, 1869.

The numbers of men present are represented by the widths of the colored zones at a rate of one millimeter for every ten thousand men; they are further written across the zones. The red designates the men who enter Russia, the black those who leave it. — The information which has served to draw up the map has been extracted from the works of M.M. Thiers, de Ségur, de Fezensac, de Chambray and the unpublished diary of Jacob, the pharmacist of the Army since October 28th. In order to better judge with the eye the diminution of the army, I have assumed that the troops of Prince Jérôme and of Marshal Davout, who had been detached at Minsk and Mogilev and have rejoined near Orsha and Vitebsk, had always marched with the army.



GRAPHIC TABLE of the temperature in degrees below zero of the Réaumur thermometer.

The Cossacks pass the frozen Neman at a gallop.

	${}^{\circ}\text{R}$	${}^{\circ}\text{C}$	${}^{\circ}\text{F}$
October 18	0	0	32
Rain October 24	-10	-13	10
December 6	-30	-38	-36
December 7	-26		
November 14	-21		
November 9	-9		
November 28	-20		
December 1	-24		

6200 BC

1500 BC

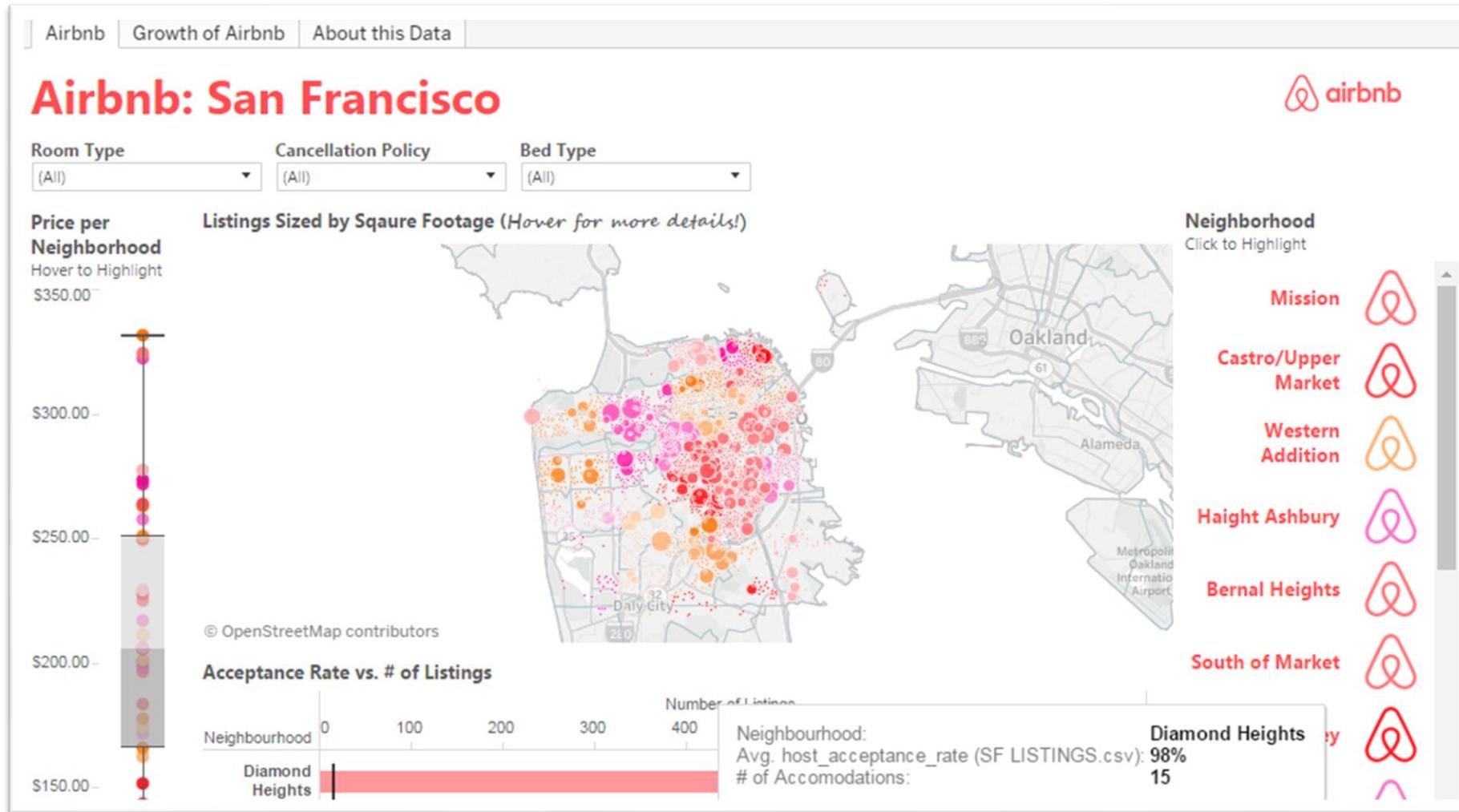
1786

1801

1855

1861

Interactive Visualization



173921017

238745328

701285647

173921017

238745328

701285647

Color

777771392

101238453

280128564

Position

173921017

238745328

Size

701285647

ແລ້ວກົມາຈົບຖ່ຽນ

ຄຸນຈະອ່ານຂ້ອຄວາມ ຕຽບນີ້ກ່ອນ

ແລ້ວກົມະນາອ່ານຂ້ອຄວາມນີ້
ຕ່ອດວຍວັນນີ້

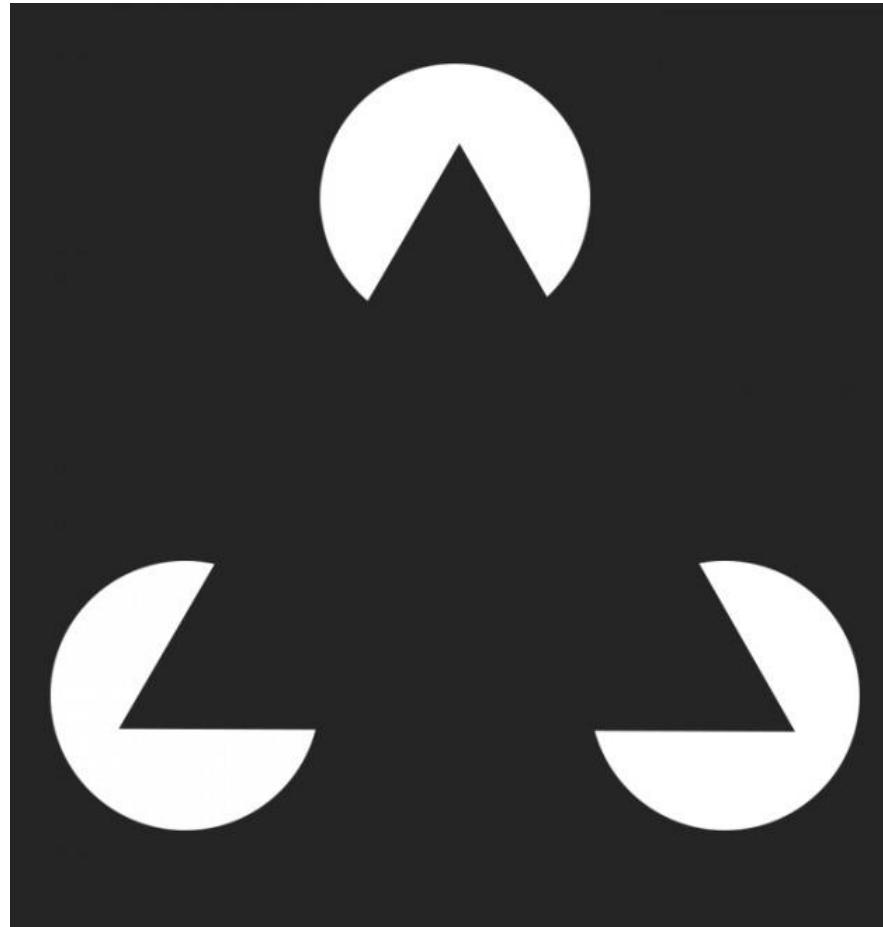
Color

Size

Enclosure

Position

Gestalt Principles of Visual Perception



Source: *The Inspired Eye*



Source: *Gizmodo*

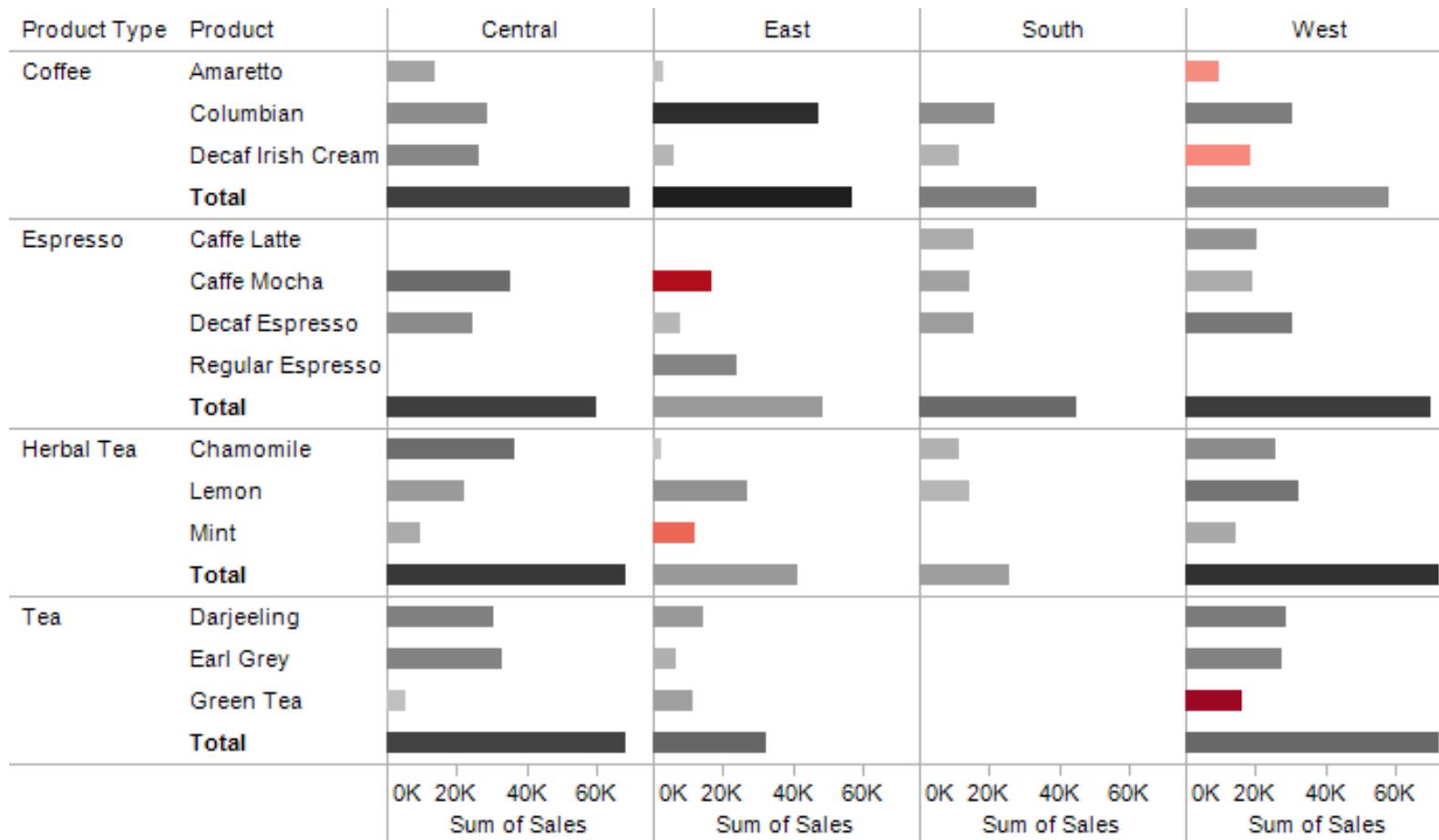
Bertin's Three Levels of Reading

Elementary: single value

Product Type	Product	Central		East		South		West	
		Sum of Profit	Sum of Sales						
Coffee	Amaretto	\$5,105	\$14,011	\$1,009	\$2,993			(\$1,225)	\$9,265
	Columbian	\$8,528	\$28,913	\$27,253	\$47,386	\$8,767	\$21,664	\$11,253	\$30,357
	Decaf Irish Cream	\$9,632	\$26,155	\$2,727	\$6,261	\$2,933	\$11,592	(\$1,305)	\$18,235
	Total	\$23,265	\$69,080	\$30,989	\$56,640	\$11,700	\$33,256	\$8,724	\$57,856
Espresso	Caffe Latte					\$3,872	\$15,442	\$7,502	\$20,458
	Caffe Mocha	\$14,640	\$35,218	(\$6,230)	\$16,646	\$5,201	\$14,163	\$4,064	\$18,876
	Decaf Espresso	\$8,860	\$24,485	\$2,410	\$7,722	\$5,930	\$15,384	\$12,302	\$30,578
	Regular Espresso			\$10,062	\$24,036				
	Total	\$23,500	\$59,703	\$6,242	\$48,405	\$15,003	\$44,989	\$23,868	\$69,911
Herbal Tea	Chamomile	\$14,434	\$36,570	\$765	\$2,194	\$3,180	\$11,186	\$8,852	\$25,632
	Lemon	\$6,251	\$21,978	\$7,901	\$27,176	\$2,593	\$14,497	\$13,120	\$32,274
	Mint	\$4,069	\$9,337	(\$2,242)	\$11,992			\$4,330	\$14,380
	Total	\$24,754	\$67,885	\$6,424	\$41,362	\$5,774	\$25,683	\$26,301	\$72,285
Tea	Darjeeling	\$10,772	\$30,289	\$6,497	\$14,096			\$11,780	\$28,769
	Earl Grey	\$10,331	\$32,881	\$3,405	\$6,505			\$10,425	\$27,387
	Green Tea	\$1,227	\$5,211	\$5,654	\$11,571			(\$7,109)	\$16,063
	Total	\$22,330	\$68,380	\$15,557	\$32,172			\$15,097	\$72,220

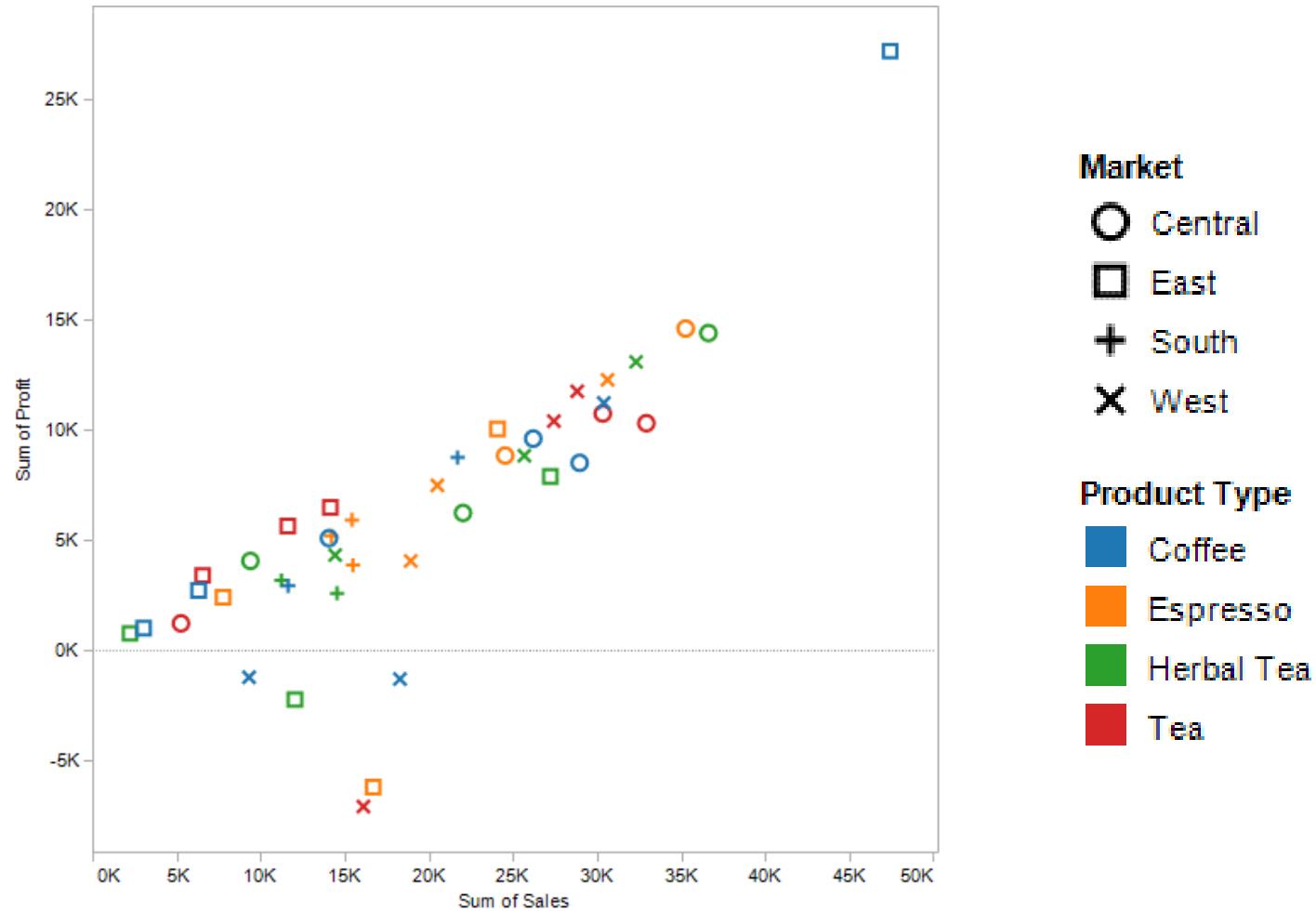
Bertin's Three Levels of Reading

Intermediate: relationships between values



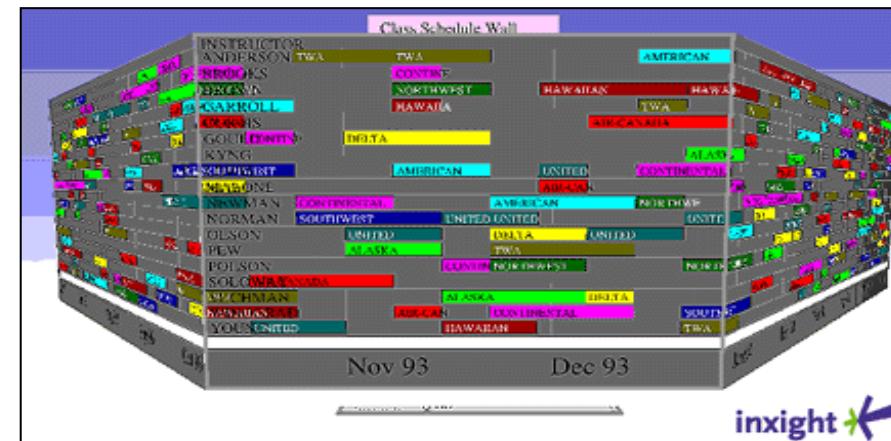
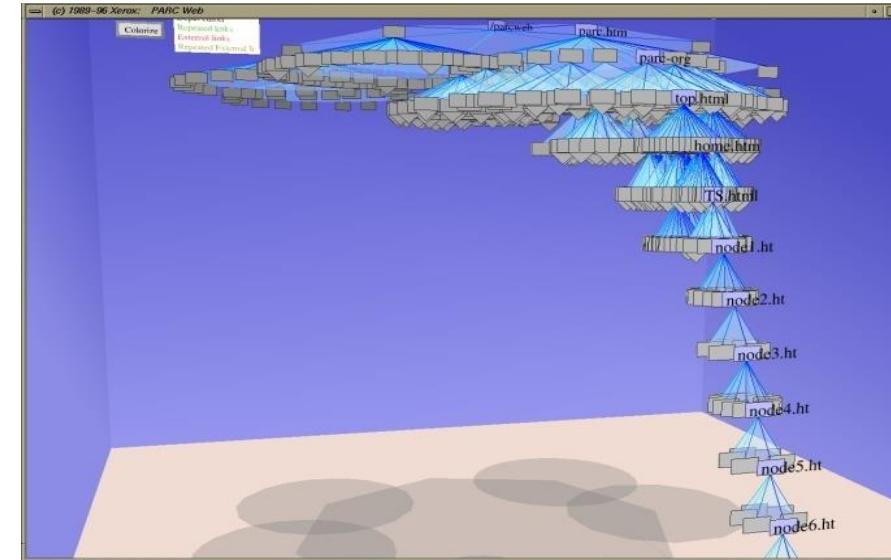
Bertin's Three Levels of Reading

Global: relationships of the whole

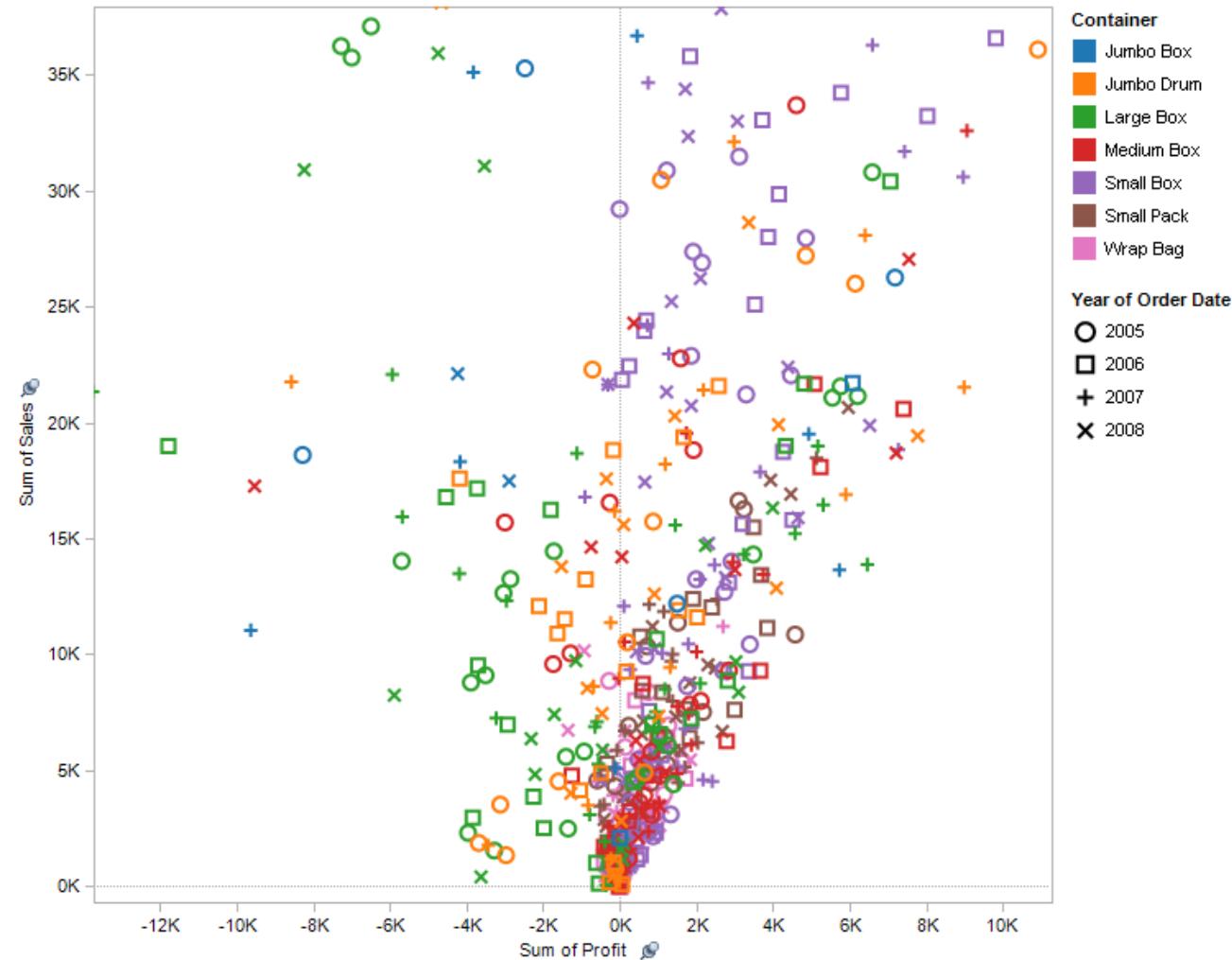


Human Perception is Limited

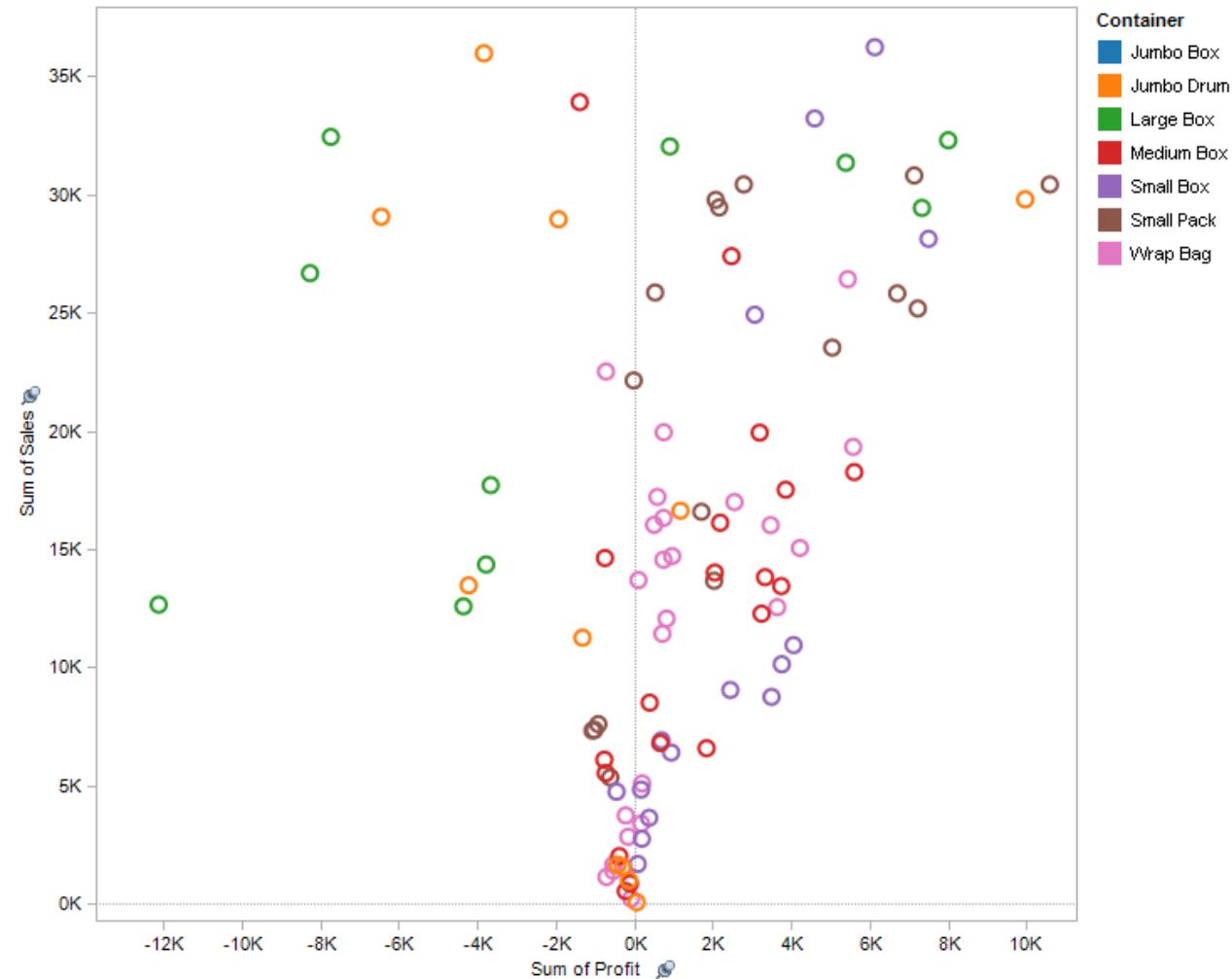
- Only adds a single dimension
- Creates occlusions
- Adds orientation complexities
- Easy to get lost
- Suggests a physical metaphor



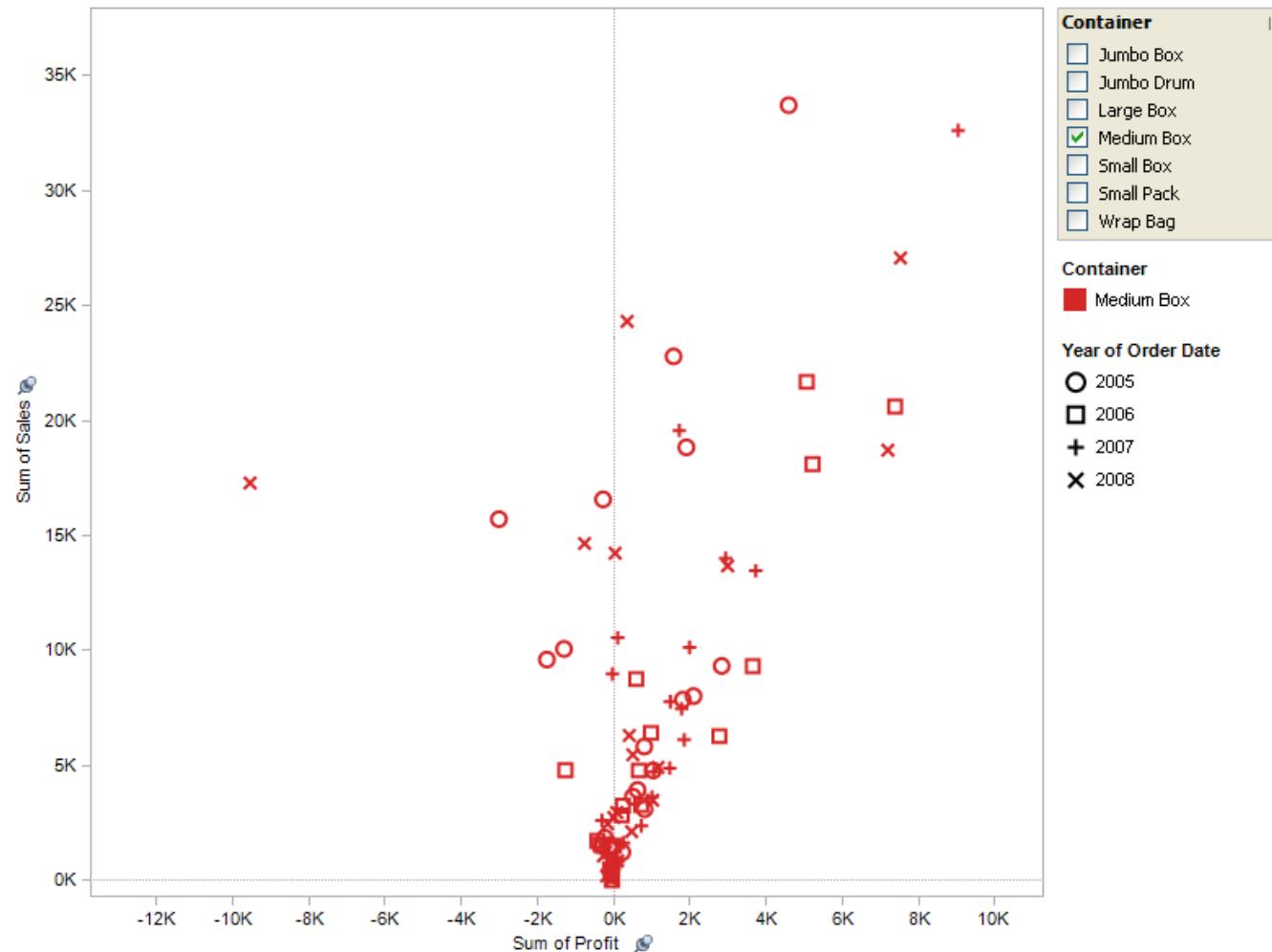
Interactivity: Too Much Data Scenario



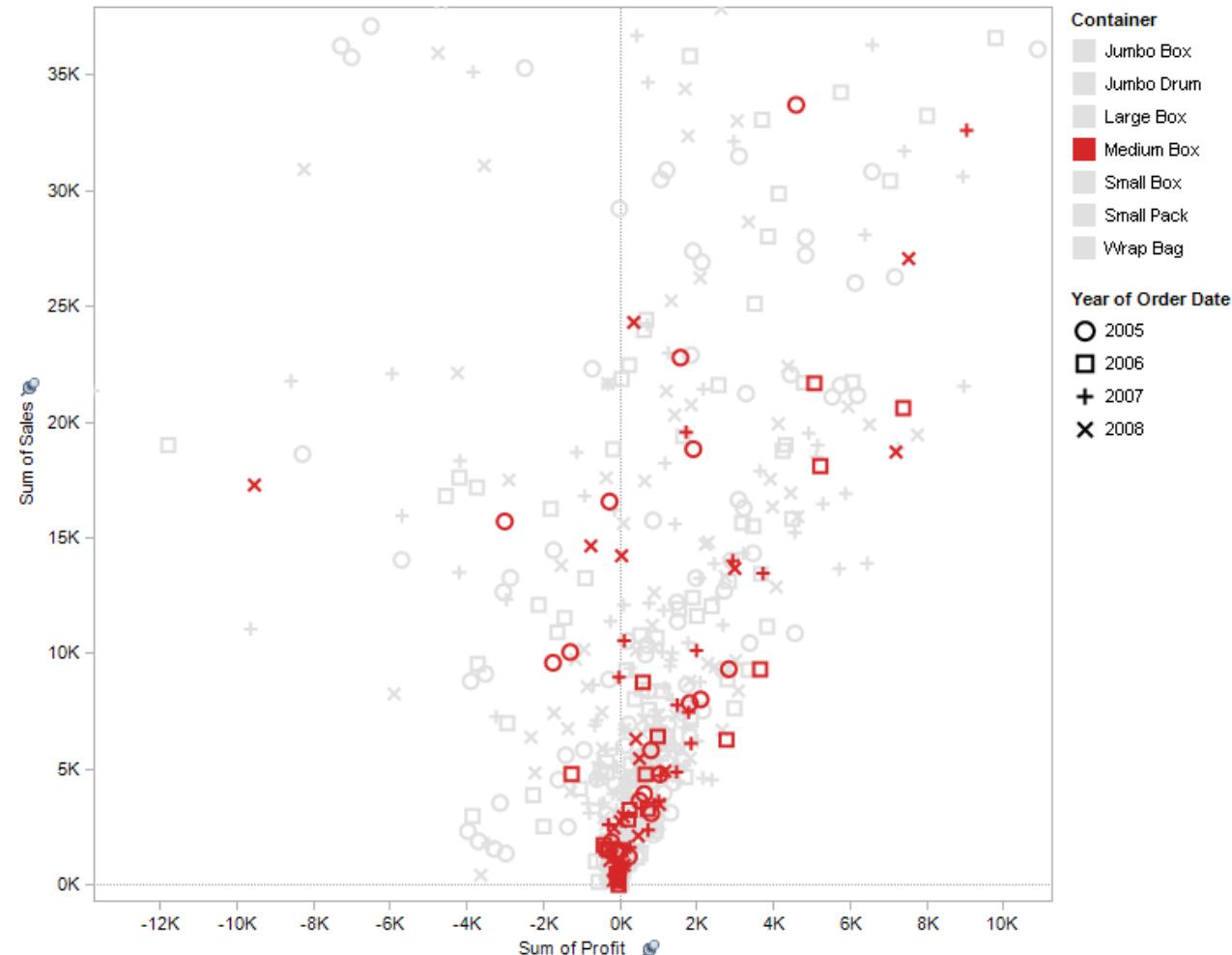
Interactivity: Aggregation



Interactivity: Filtering



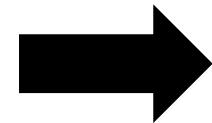
Interactivity: Brushing



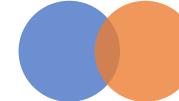
What do you want to tell from your data?



Story



1.Comparison



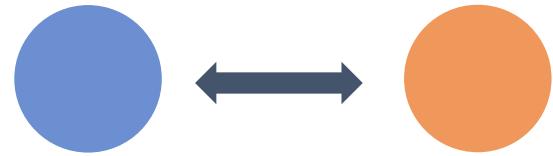
2.Relationship



3.Composition



4.Distribution



1. Comparison

Bar chart

Line chart

Bullet chart

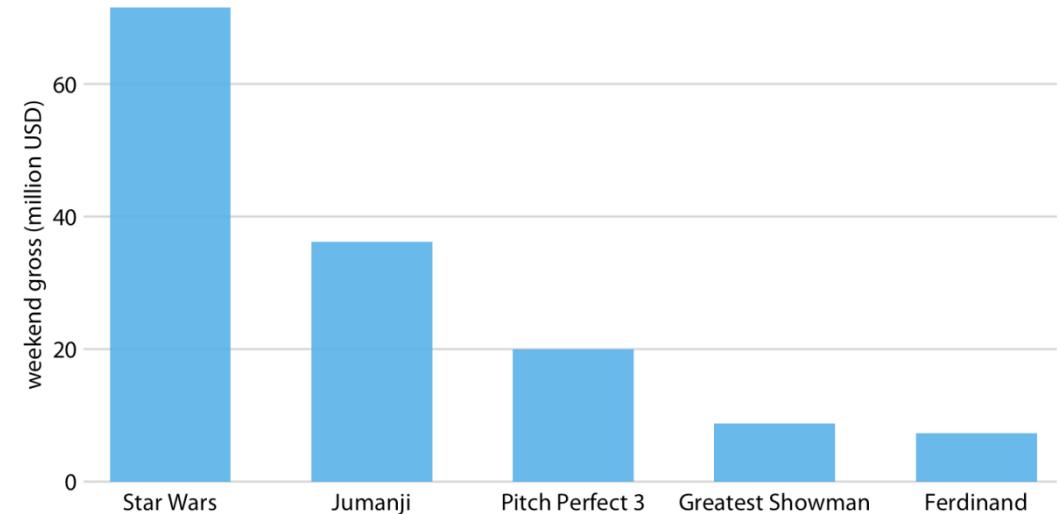
Bar chart

When to use

- Comparing data across categories

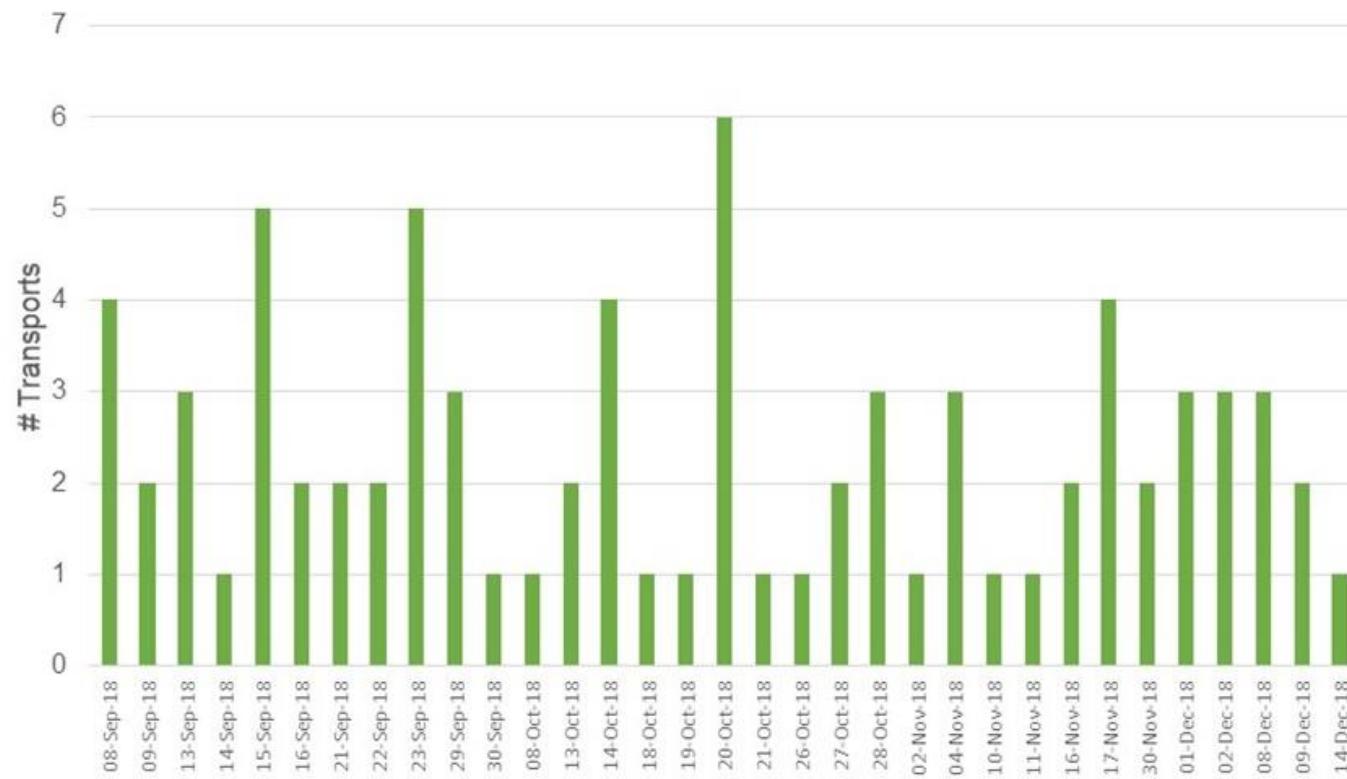
Possible extension

- Include multiple bar charts on a dashboard
- Add color to bars for more impact
- Use stacked bars or side-by-side bars
- Combine bar charts with maps
- Put bars on both sides of an axis



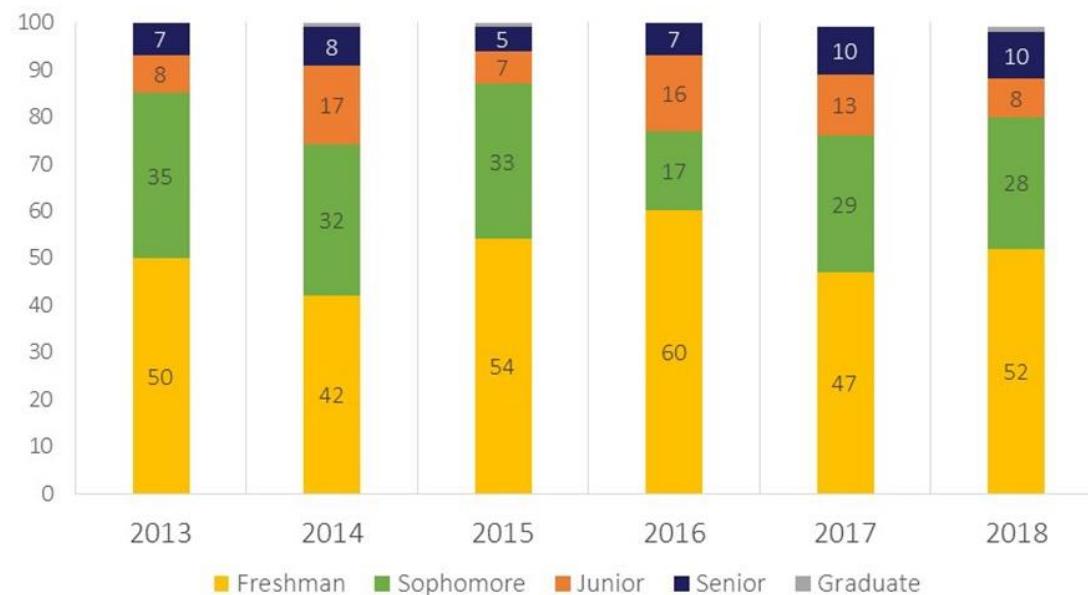
Bar chart

Transport Incident Dates

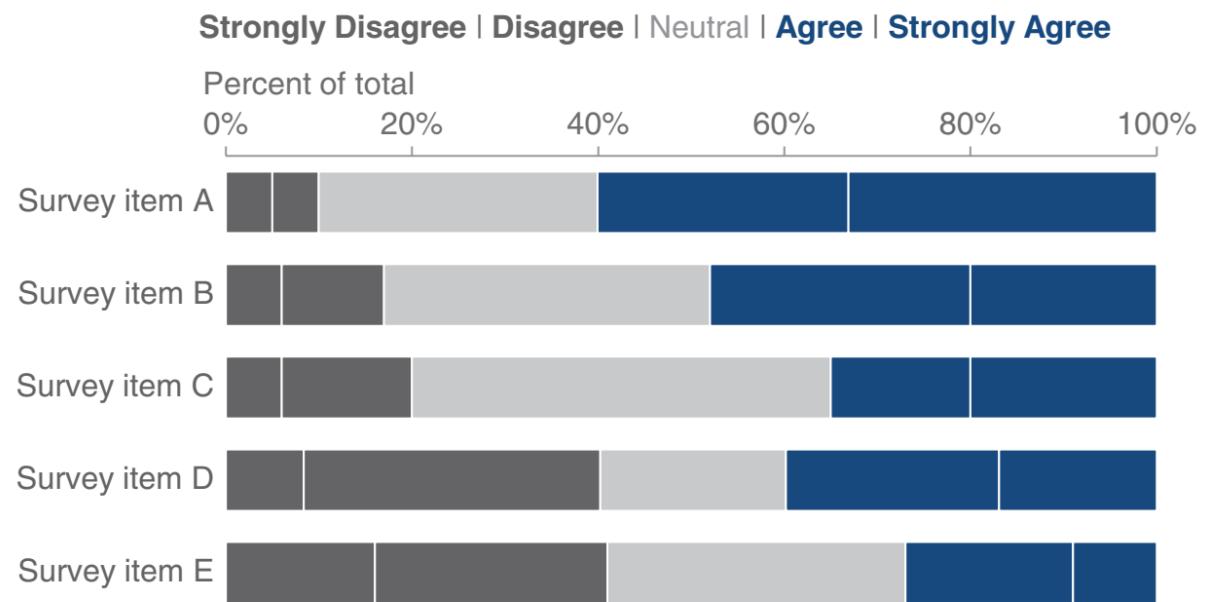


Bar chart (stacked)

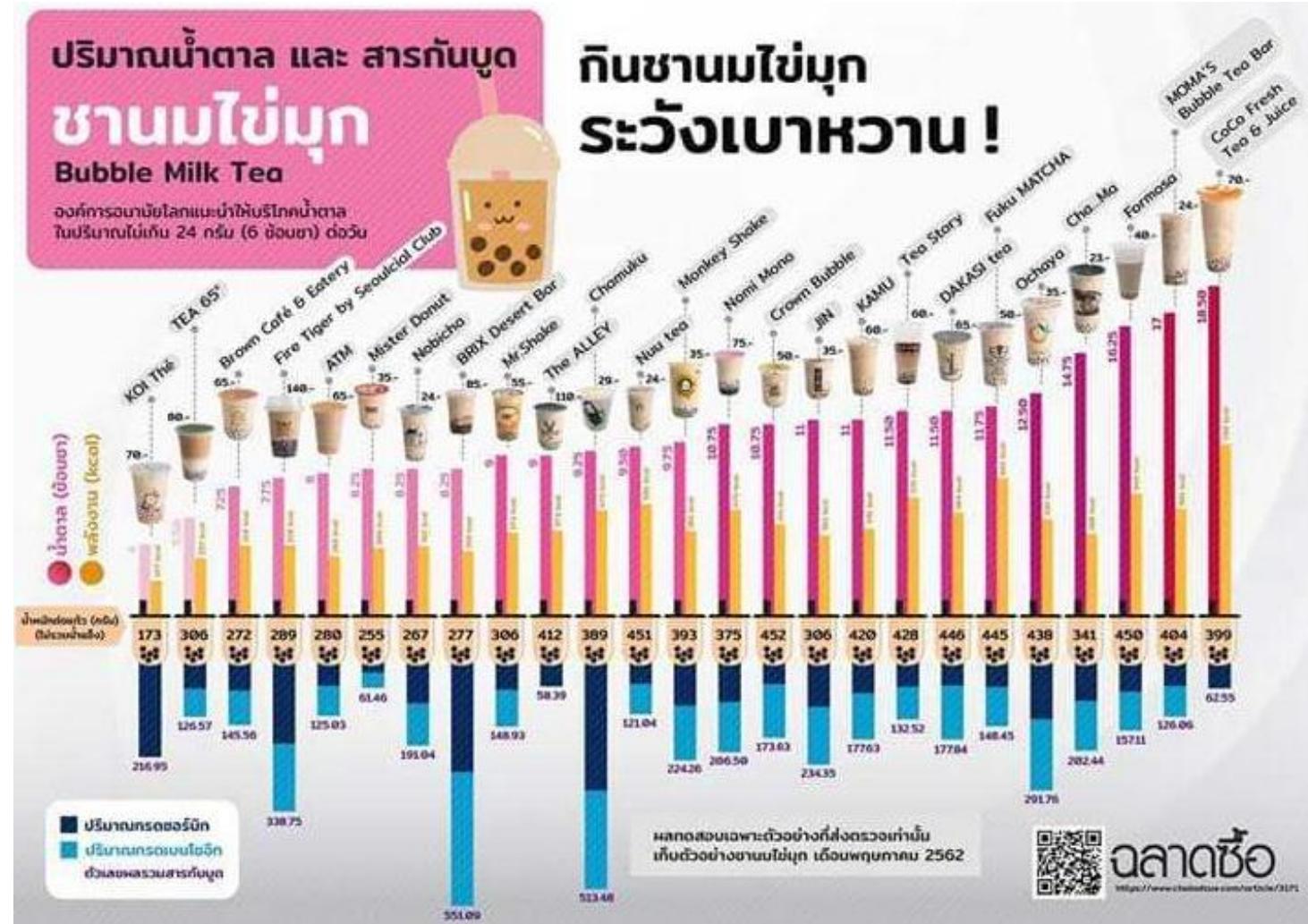
Class Year Distribution



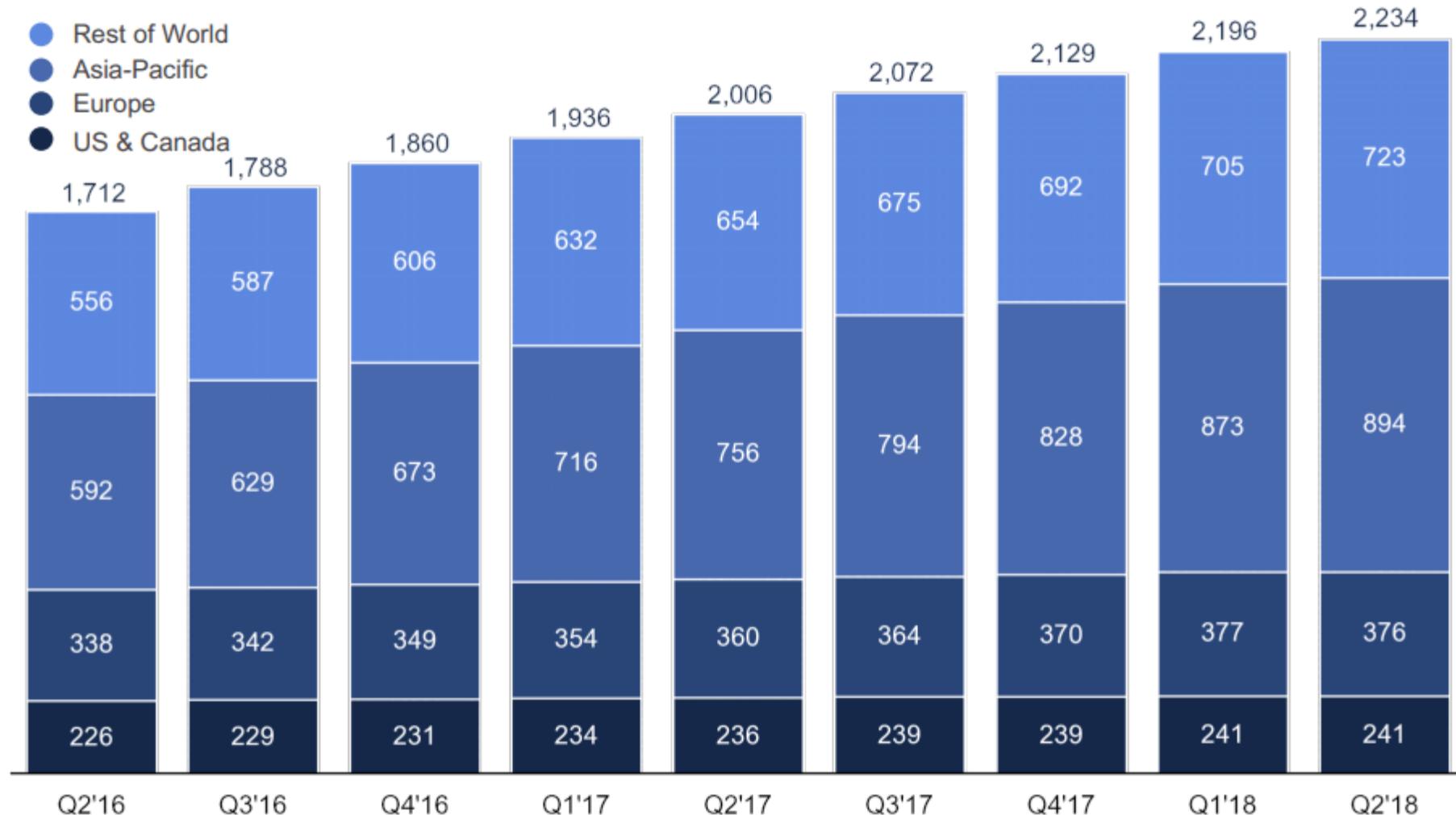
Survey results



Bar chart (side by side)



FACEBOOK STATISTICS 2018

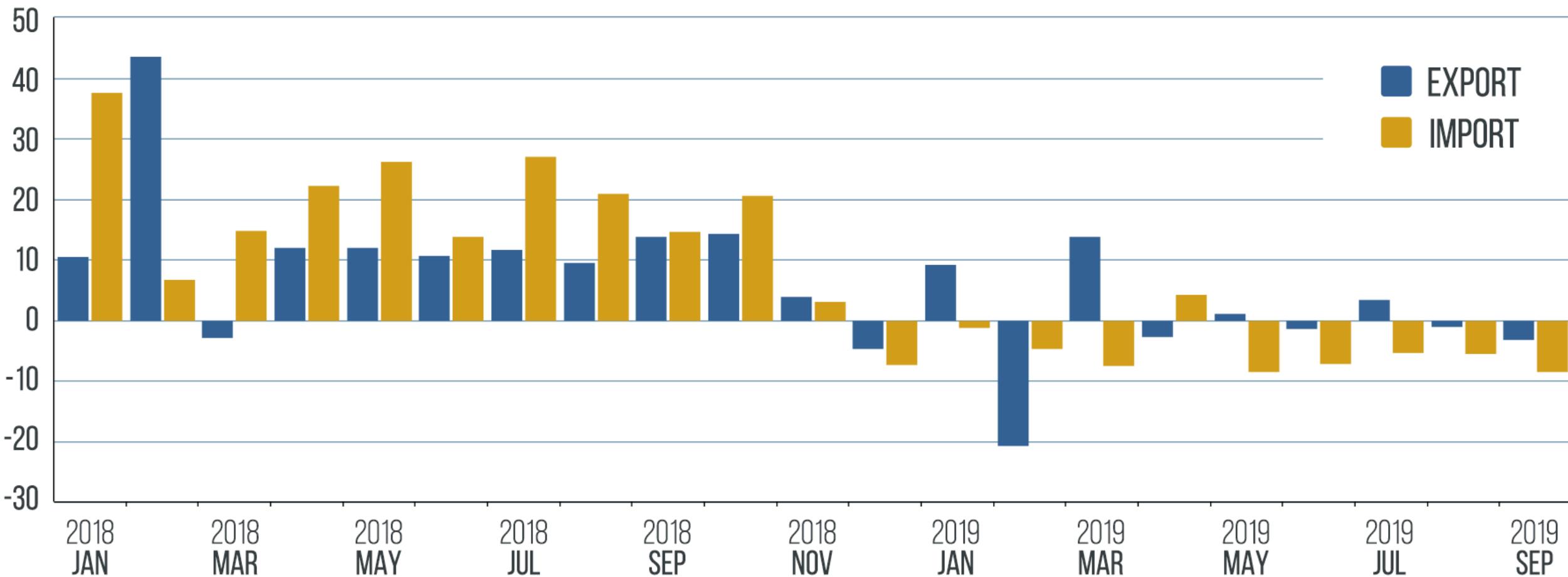


Monthly active users (MAUs), in millions

China: Export and Import Trends, Jan 2018 - Sept 2019

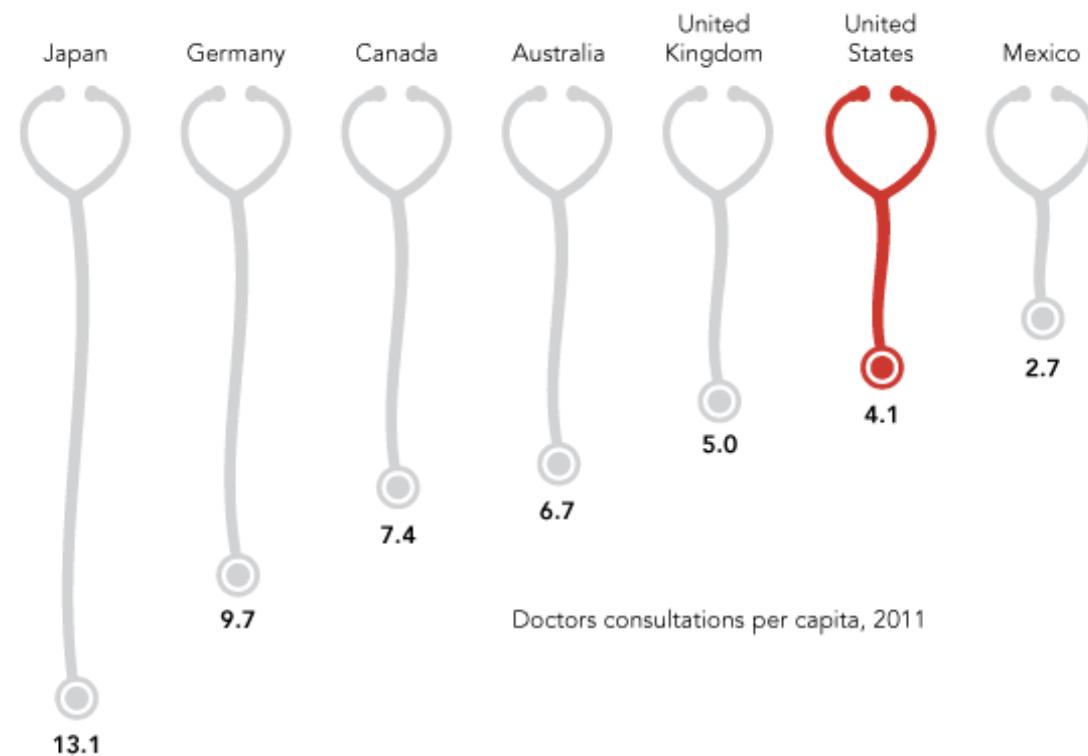
Data Driven

Y-o-Y % Change



Bar chart doesn't have to be boring

Despite high spending, Americans don't go to the doctor very frequently.

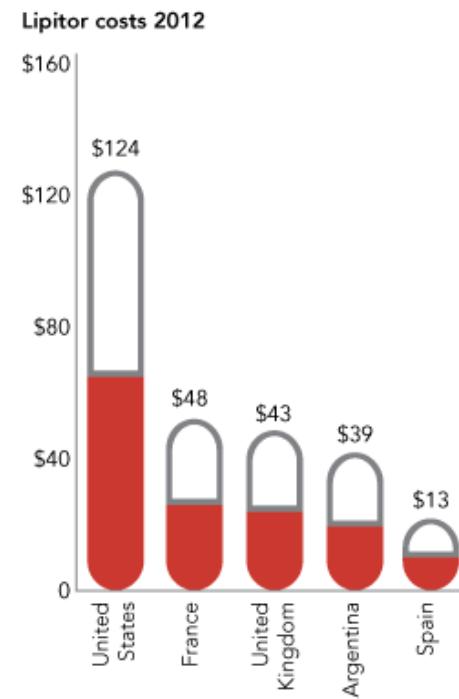
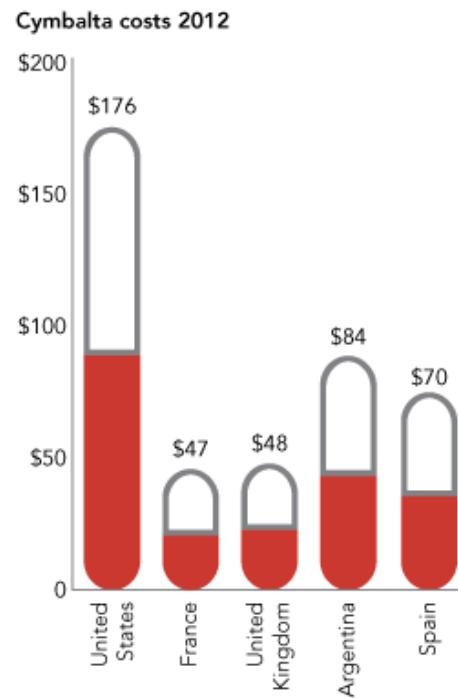


Notes: Data is from 2011 or nearest year.
Source: OECD Health Data 2013

THE HUFFINGTON POST

Bar chart doesn't have to be boring

Commonly prescribed drugs are very expensive in the U.S.



Source: International Federation of Health Plans, 2012

THE HUFFINGTON POST

- **Missed opportunity (and can be misleading)**
- **Not a stacked bar chart but could be.**

 2561
 2562

สถิติอุบัติเหตุ

วันอันตราย

สงกรานต์ 2561-2562

(วันที่ 11-17 เมษายน)



อุบัติเหตุรวม

3,724 ครั้ง

3,338 ครั้ง



บาดเจ็บรวม

3,897 คน

3,442 คน



ผู้เสียชีวิต

418 ราย

386 ราย

จังหวัดที่มีผู้เสียชีวิตมากที่สุด
ลพบุรี-อุดรธานี 15 ราย



สาเหตุการเกิดอุบัติเหตุสูงสุด

เมื่อแล้วขับ

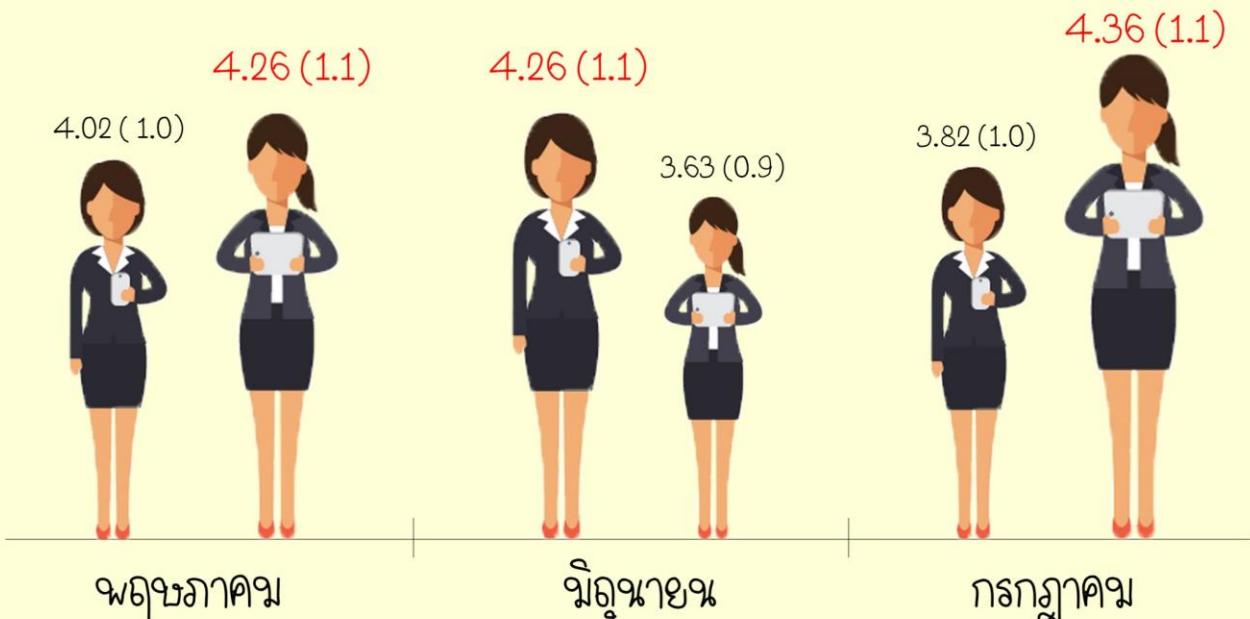


ที่มา : ศูนย์อำนวยความปลอดภัยทางถนน



ເຊື້ອງນາຍົກລົງແລະອໍ້ມະກາງວິຊາ

ພ.ສ. 2561 ແລະ 2562



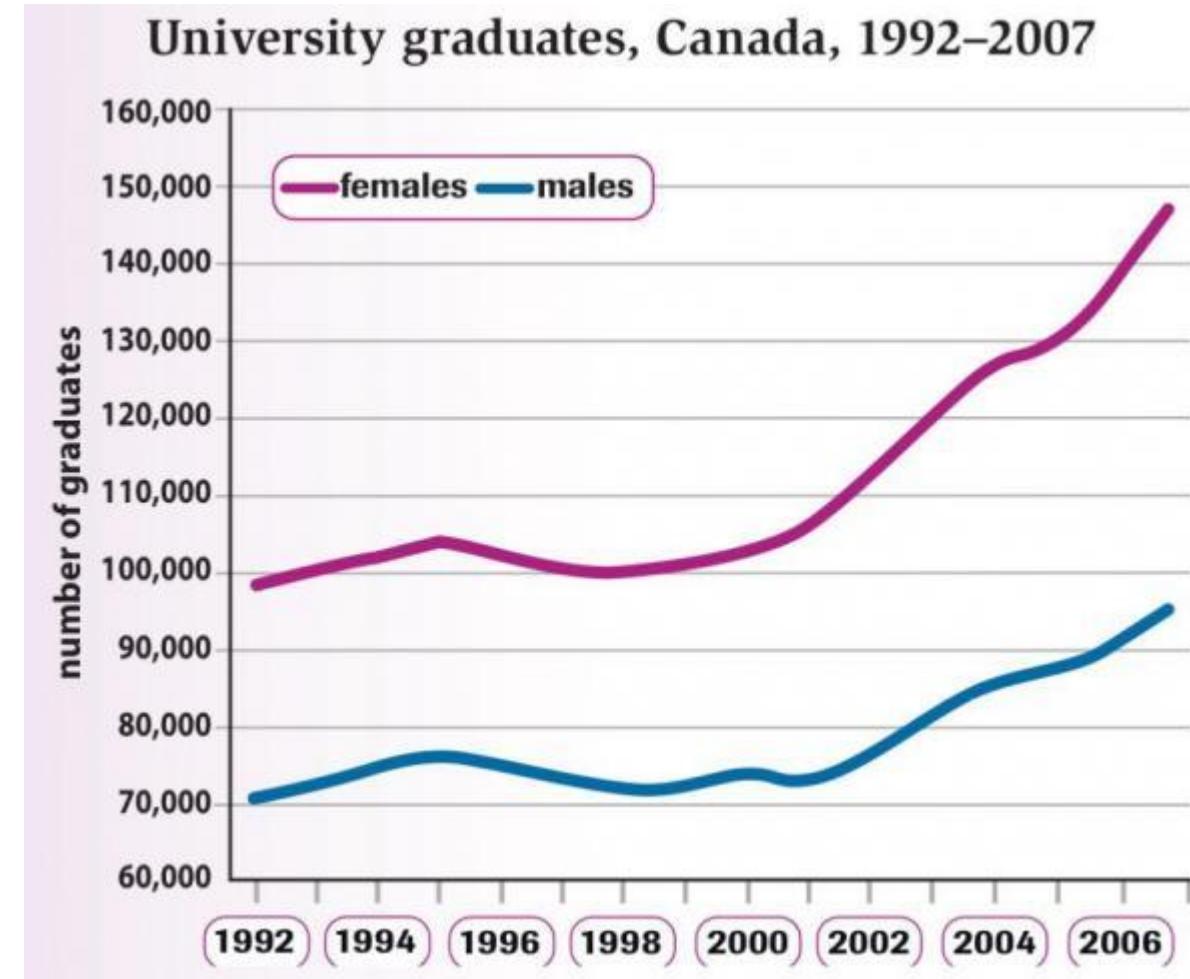
Line chart

When to use

- Viewing trends in data over time

Possible extension

- Combine a line graph with bar charts
- Shade the area under lines.



Line chart

When to use

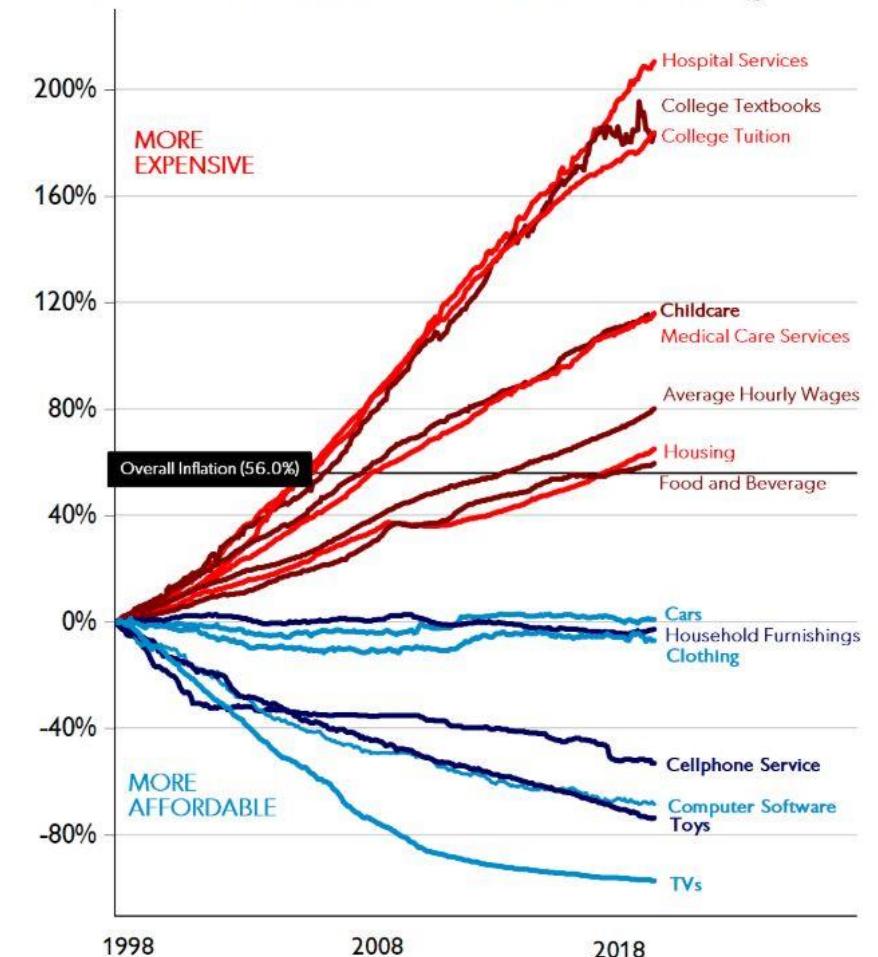
- Viewing trends in data over time

Possible extension

- Use different types/colors of lines to signify different data/meanings
- Combine a line graph with bar charts
- Shade the area under lines (Area chart).

Price Changes (January 1998 to December 2018)

Selected US Consumer Goods and Services, Wages

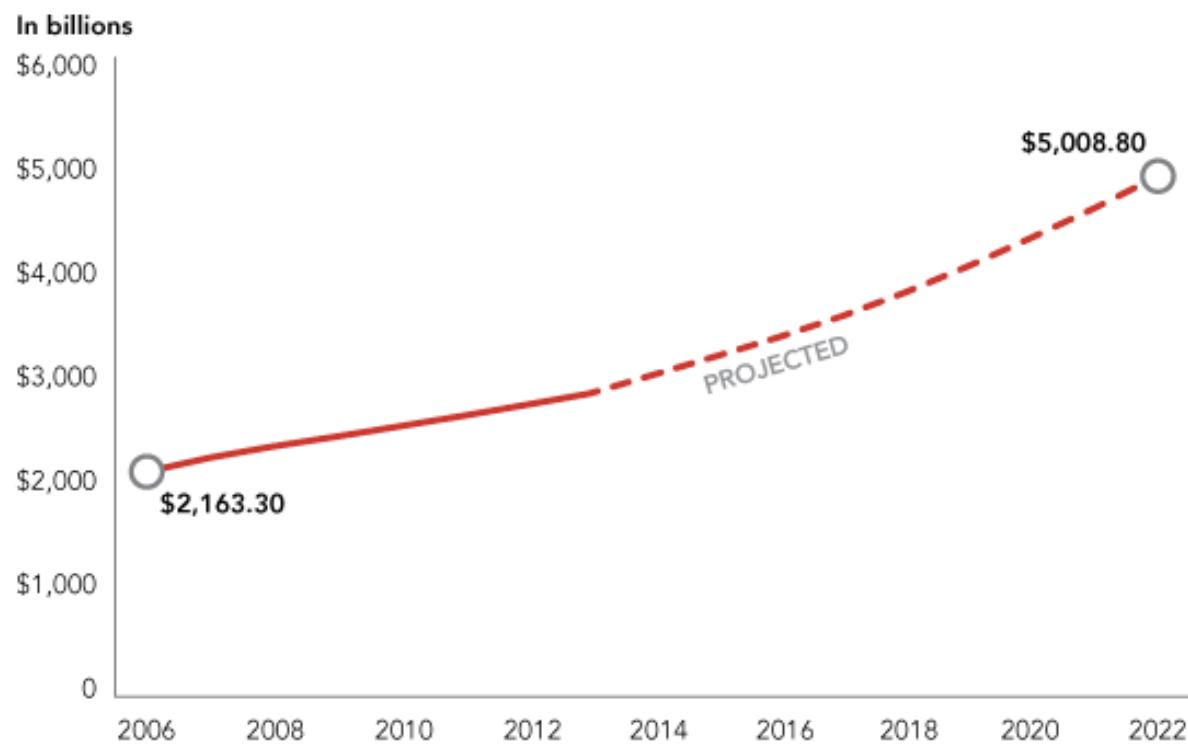


Source: BLS

Carpe Diem AEI

Line chart: different type of lines

Health care spending is projected to nearly double in the next decade.

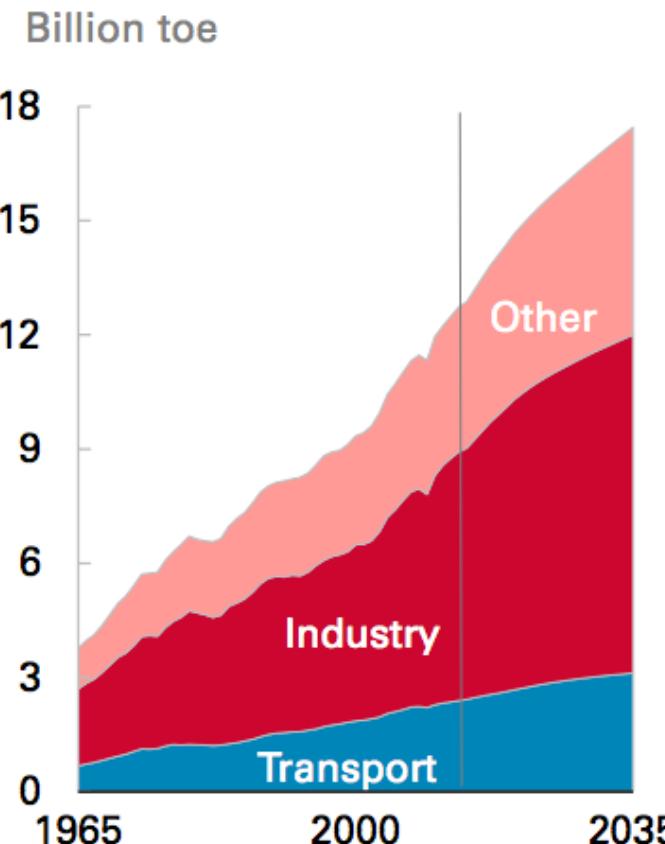


Notes: The health spending projections were based on the National Health Expenditures released in January 2013. The projections include impacts from the Affordable Care Act. Numbers may not add to totals because of rounding.
Source: Centers for Medicare & Medicaid Services, Office of the Actuary

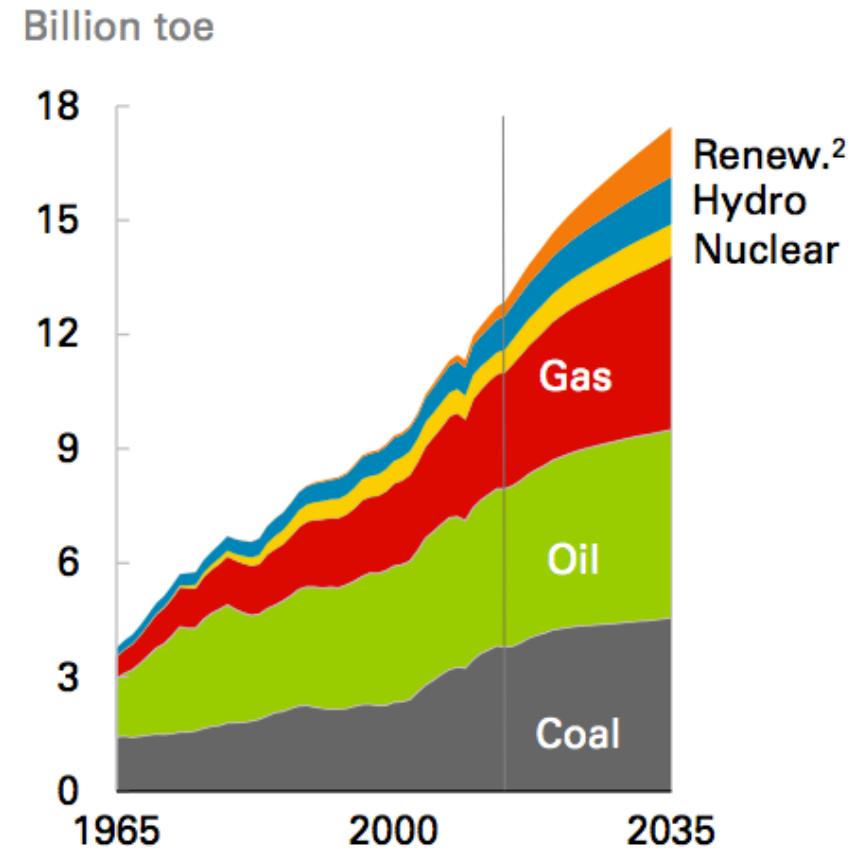
THE HUFFINGTON POST

Line chart (area)

Consumption by final sector¹



Consumption by fuel

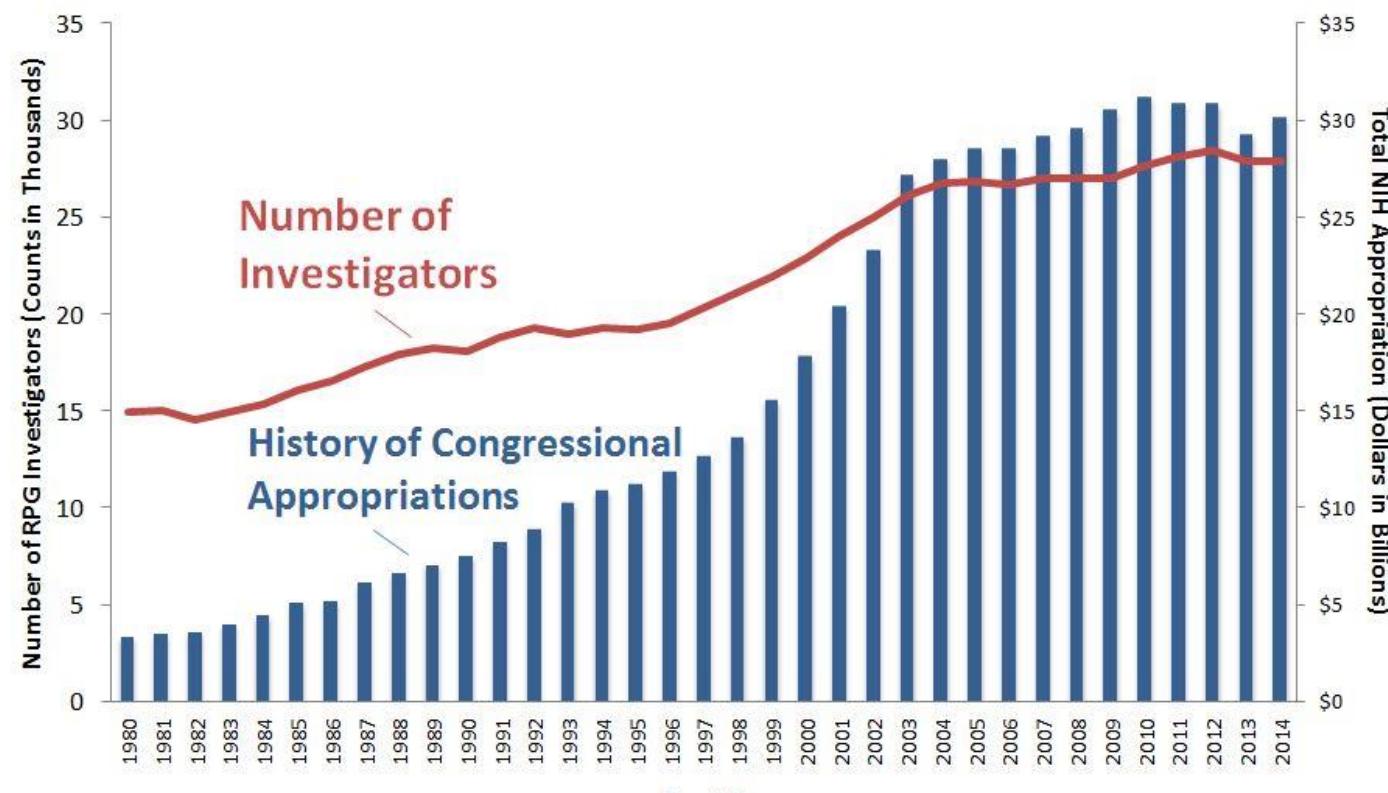


¹Primary fuels in power allocated according to final sector electricity consumption

²Includes biofuels

Bar chart + Line chart

Number of Principal Investigators* Supported on NIH Research Project Grants (RPGs) and History of Congressional Appropriations

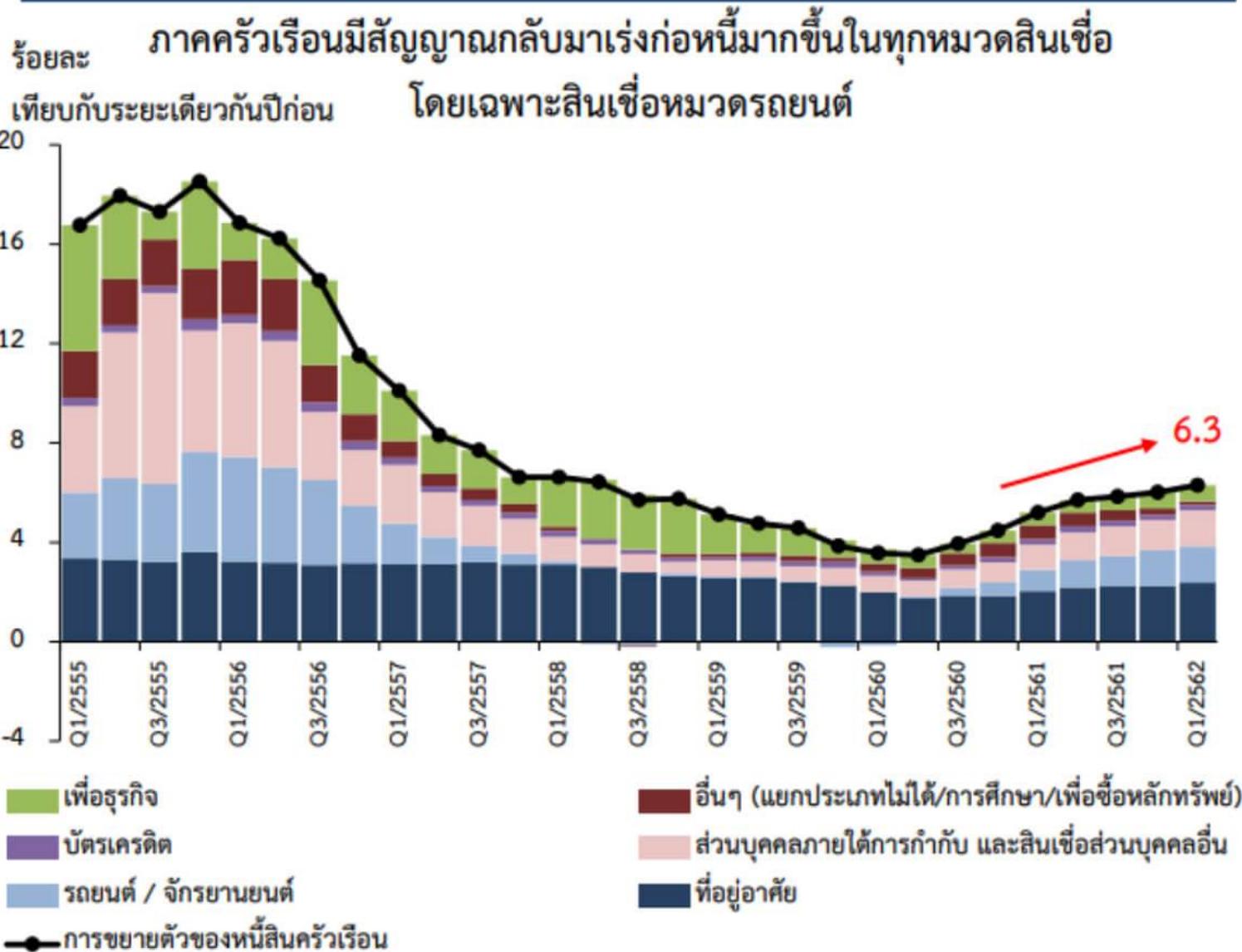


NIH Rock Talk Blog: <http://nexus.od.nih.gov/all/category/blog/>
NIH RePORT: http://report.nih.gov/special_reports_and_current_issues/index.aspx

Fiscal Year

*Includes contact and multiple principal investigators. Excludes awards made with American Recovery and Reinvestment Act funds.

แหล่งที่มาของการขยายตัวหนี้ครัวเรือนแยกตามวัตถุประสงค์



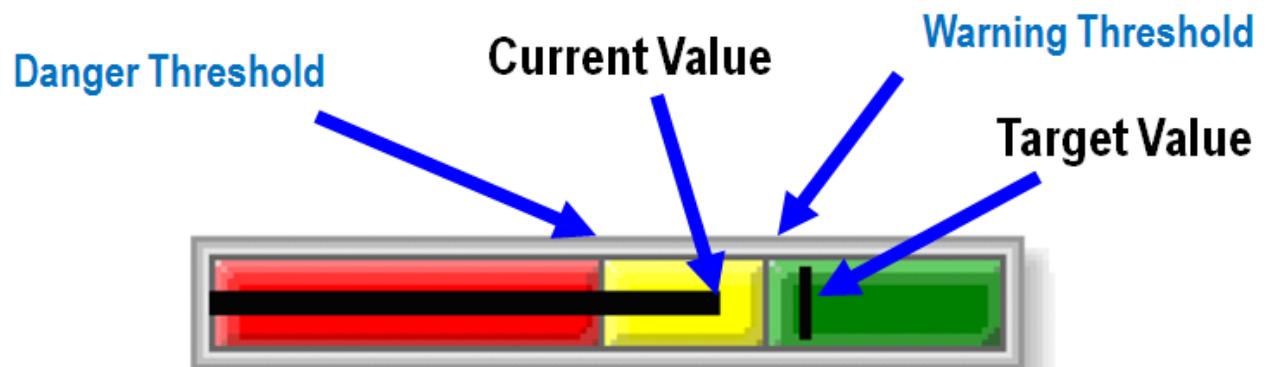
Bullet chart

When to use

- Evaluating performance of a metric against a goal

Possible extension

- Use color to illustrate achievement thresholds

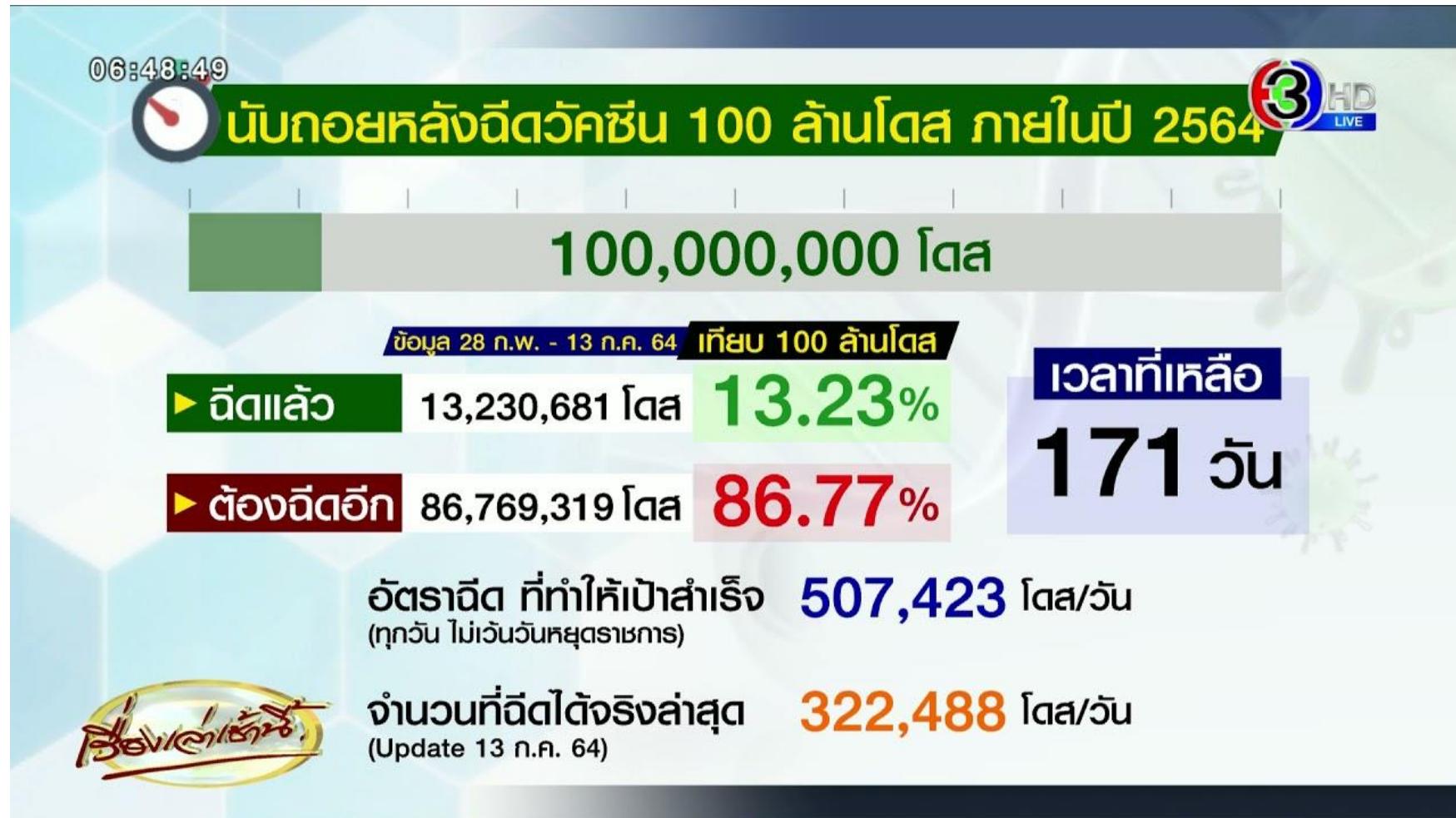


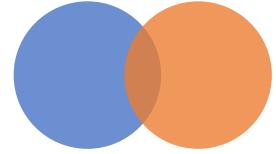
Bullet chart



Have you hit your target?

Bullet chart





2. Relationship

Scatter plot

Map

Bubble chart

Heat map

Crosstab / Highlight table

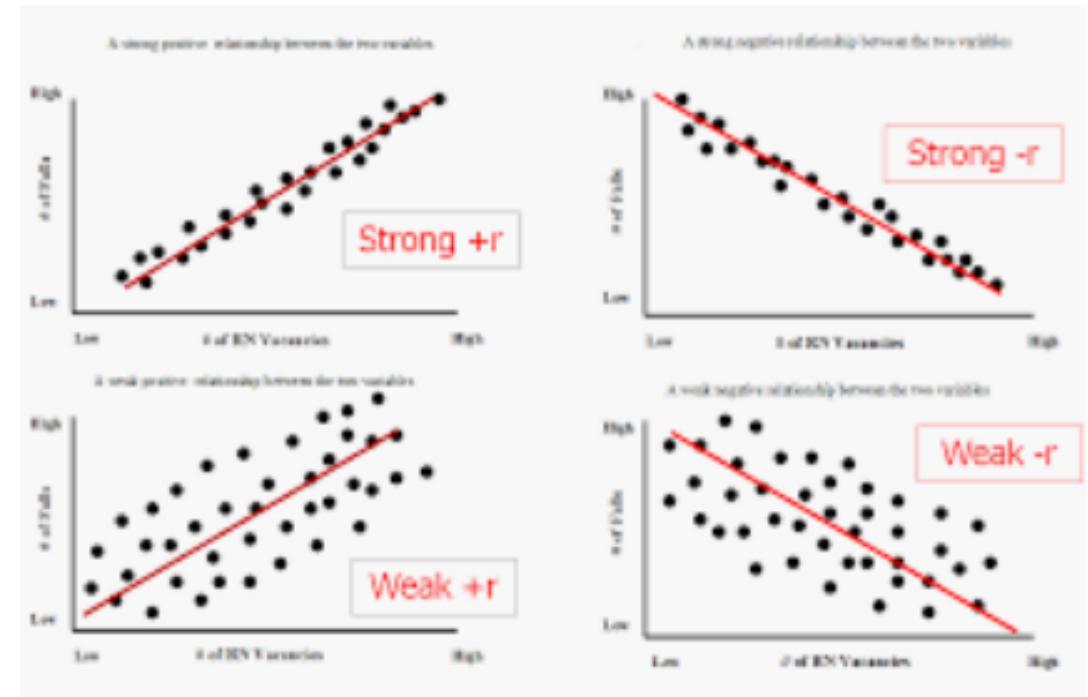
Scatter plot

When to use

- Investigating the relationship between different variables

Possible extension

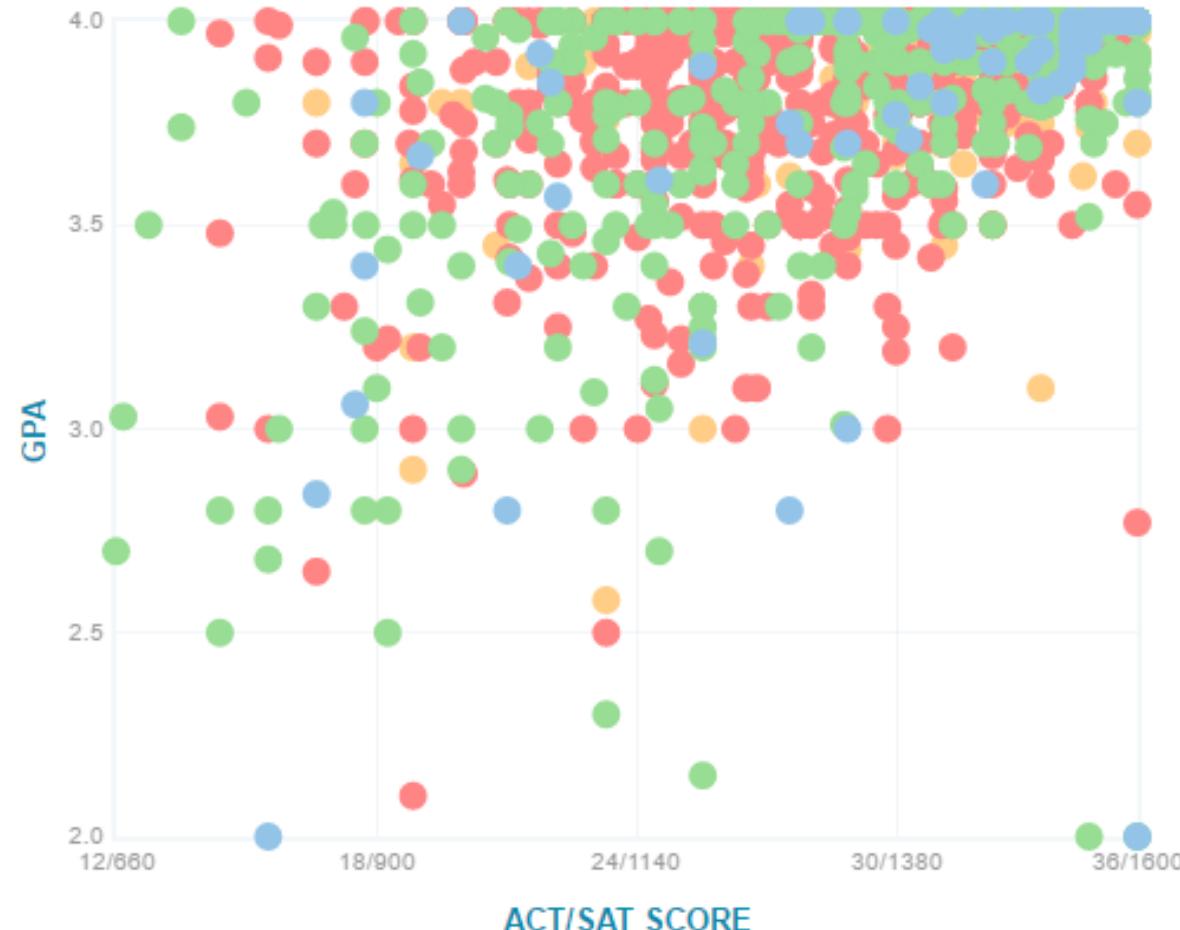
- Incorporate filters
- Use informative mark types



Scatter plot

Harvard University Admissions Scattergram

(BASED ON HISTORICAL SELF-REPORTED STUDENT DATA)

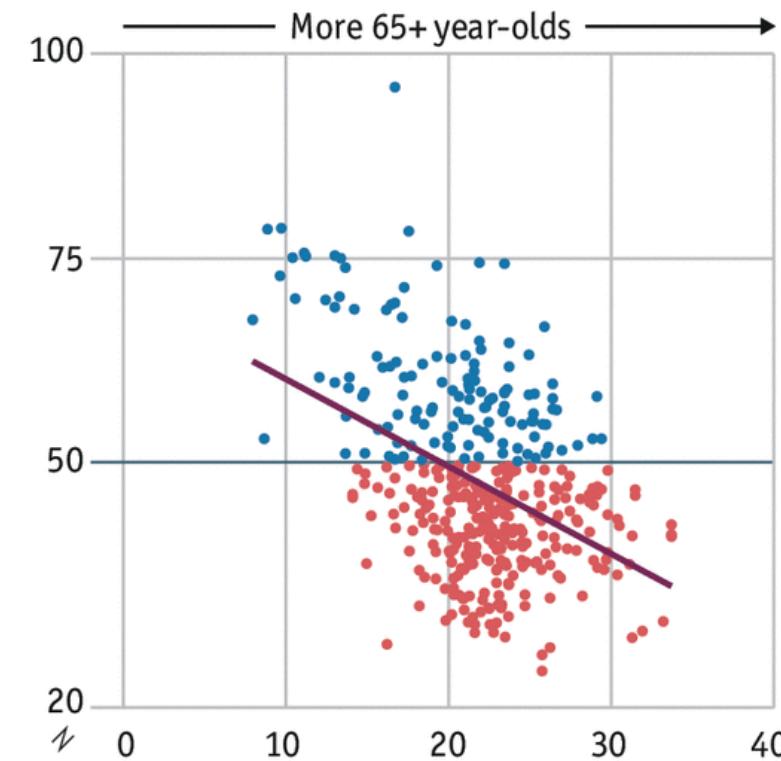
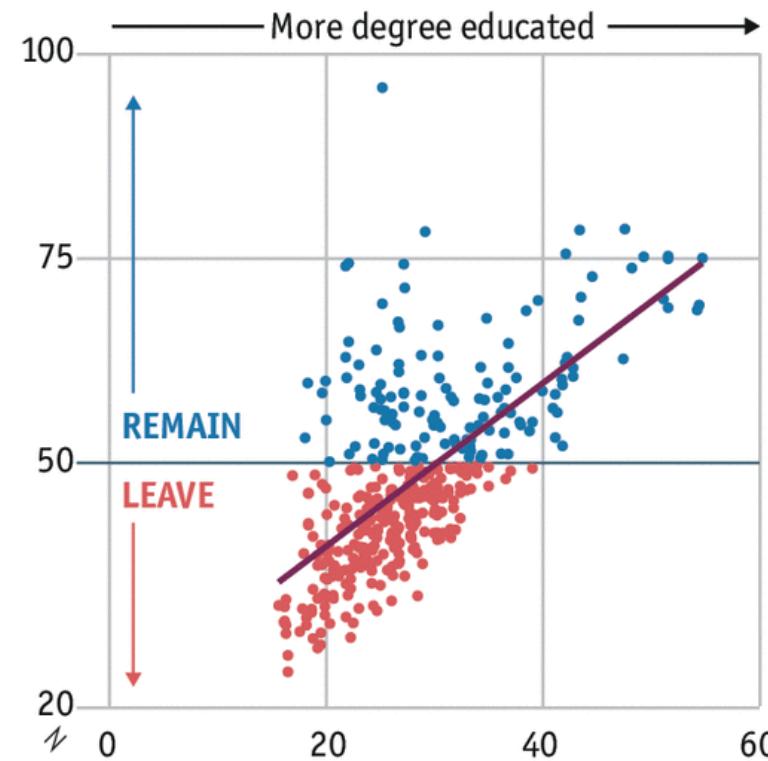


- **Green** - a successful applicant that chose to attend Harvard,
- **Blue** - an admitted applicant that chose not to attend,
- **Red** - a rejected applicant.

Scatter plot

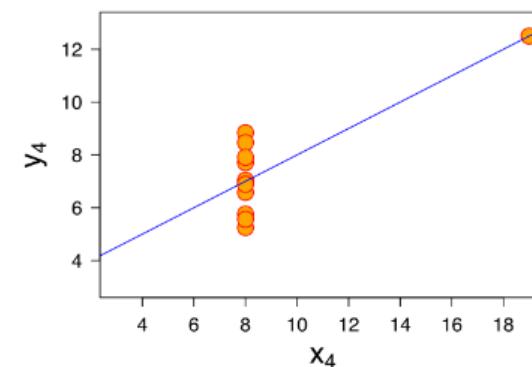
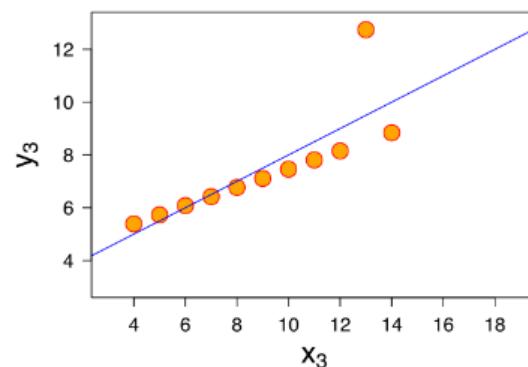
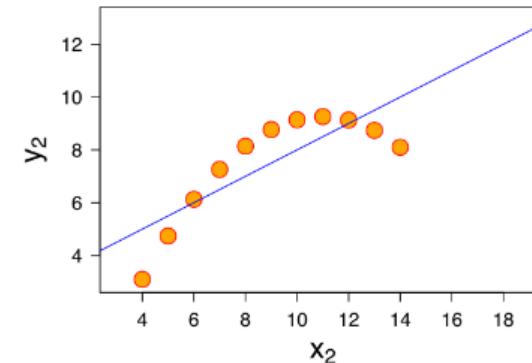
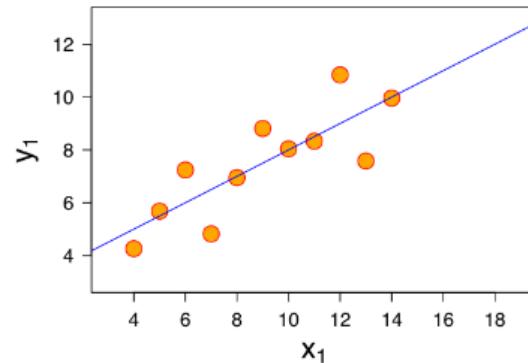
EU referendum results by demographics

Remain vote % by counting area



Sources: BBC; 2011 Census, UK Data Service

Scatter plot: Anscombe's quartet



- All 4 sets of data shares many common statistics properties: means, variance, correlations, regression line, etc.
- the importance of graphing data.
- the effect of outliers and other influential observations on statistical properties.

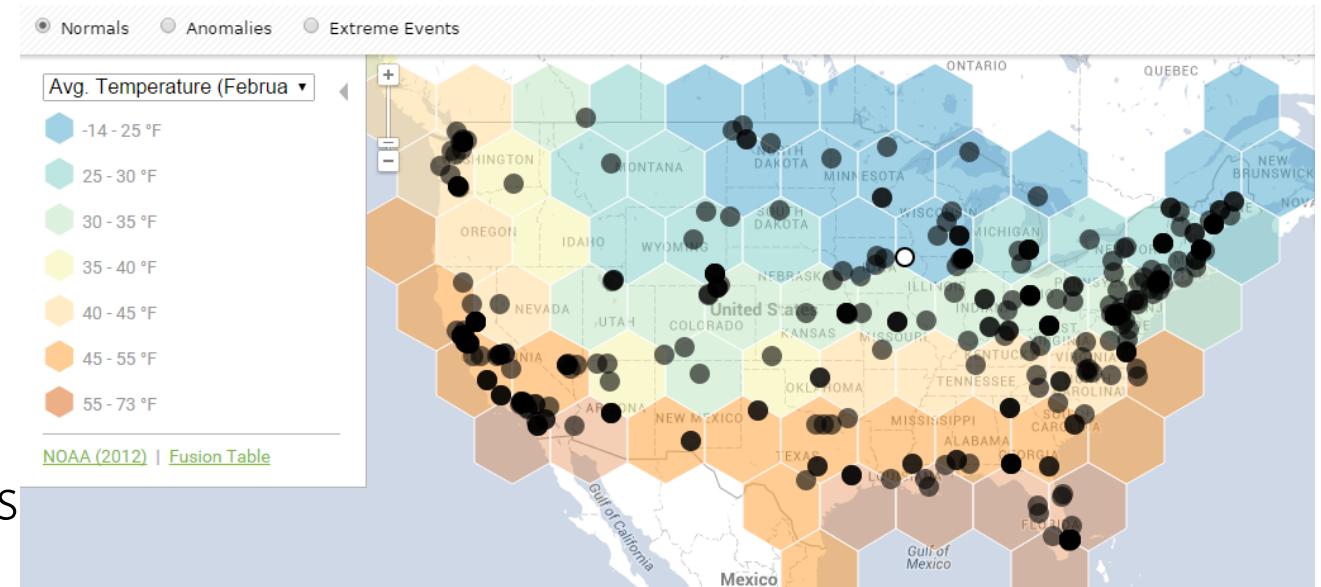
Map + Scatter plot

When to use

- Showing geocoded data

Possible extension

- Use maps as a filter for other types of charts, graphs, and tables



Climate Change news report

Bubble chart

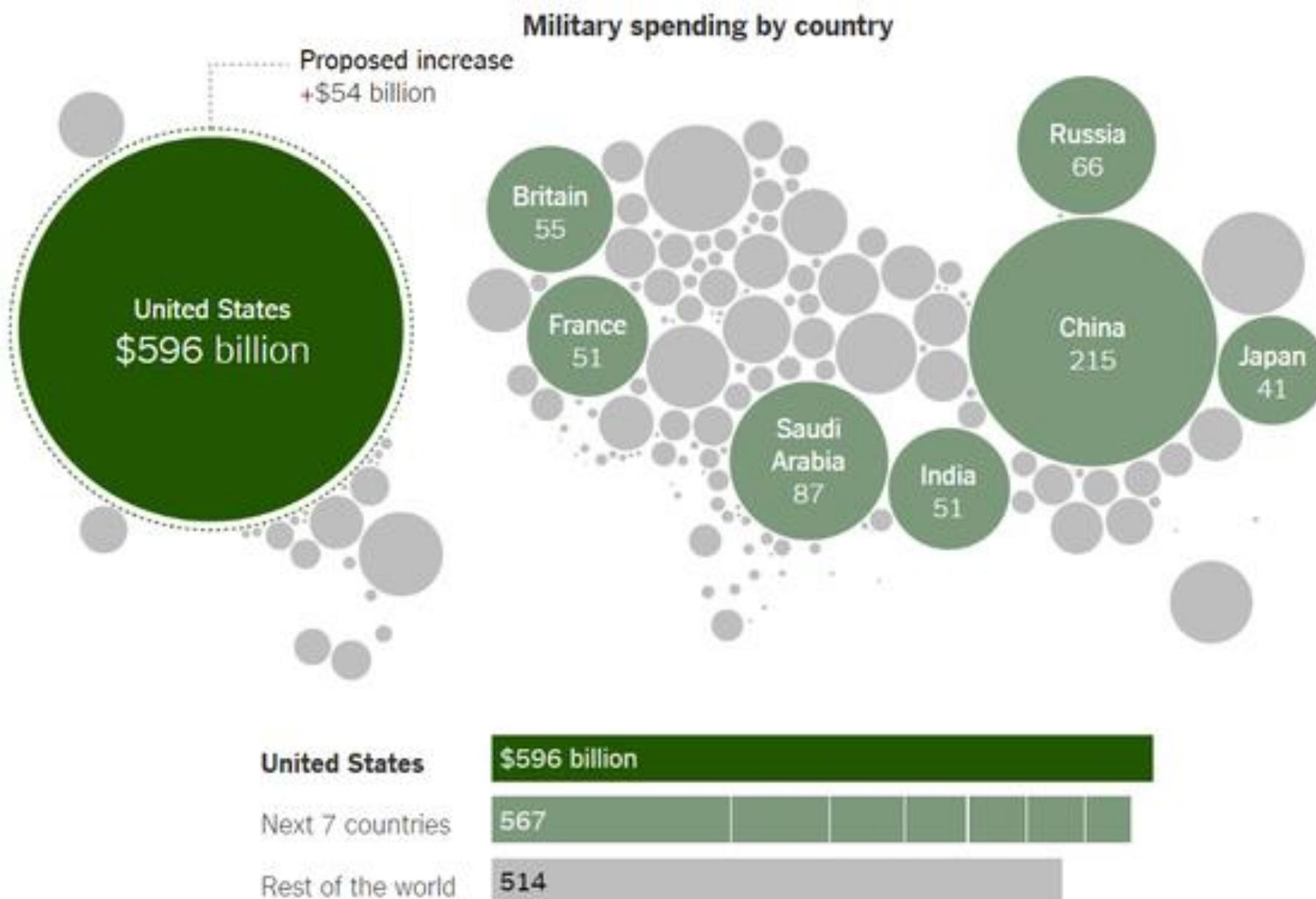
When to use

- Showing the concentration of data along two axes

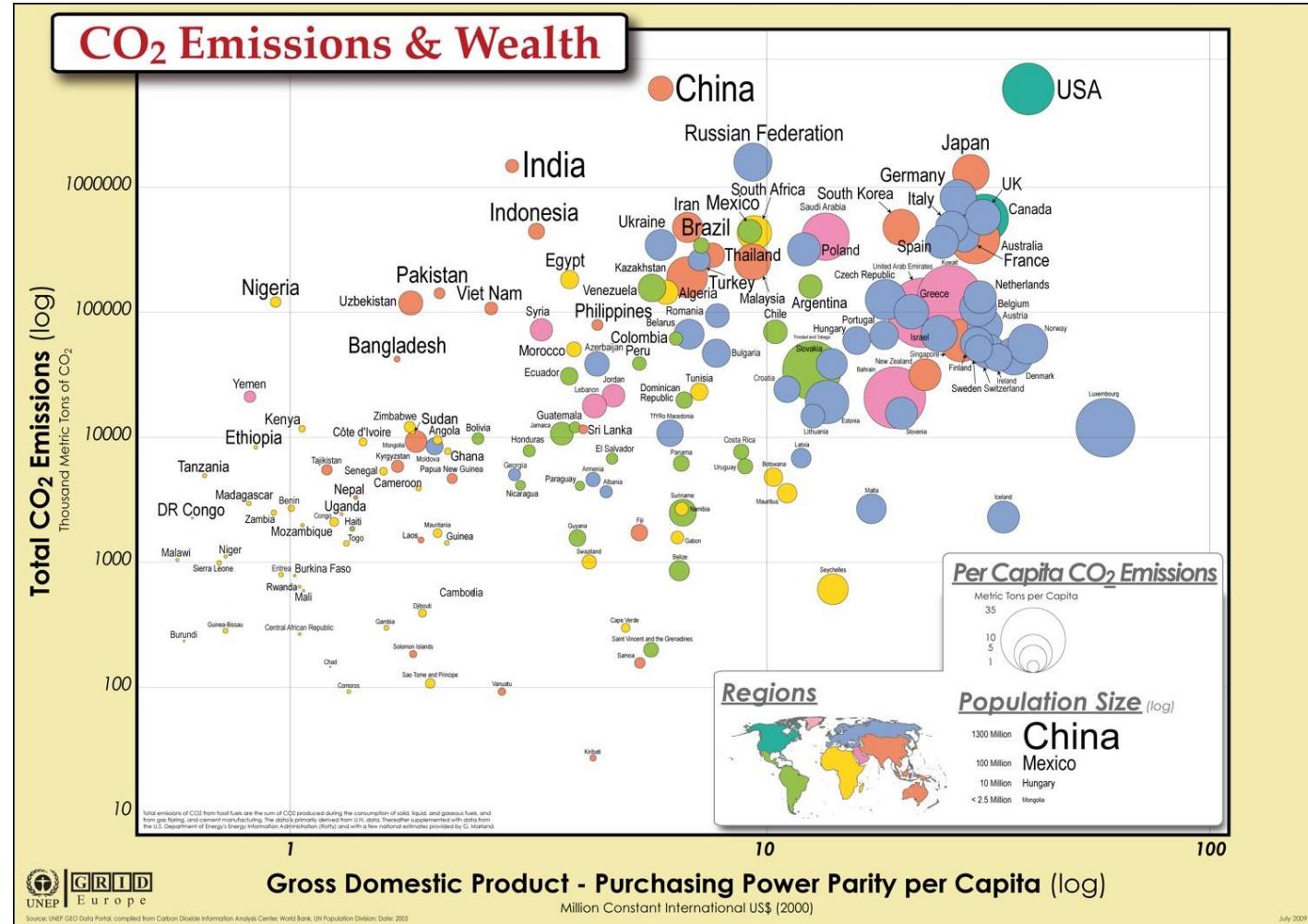
Possible extension

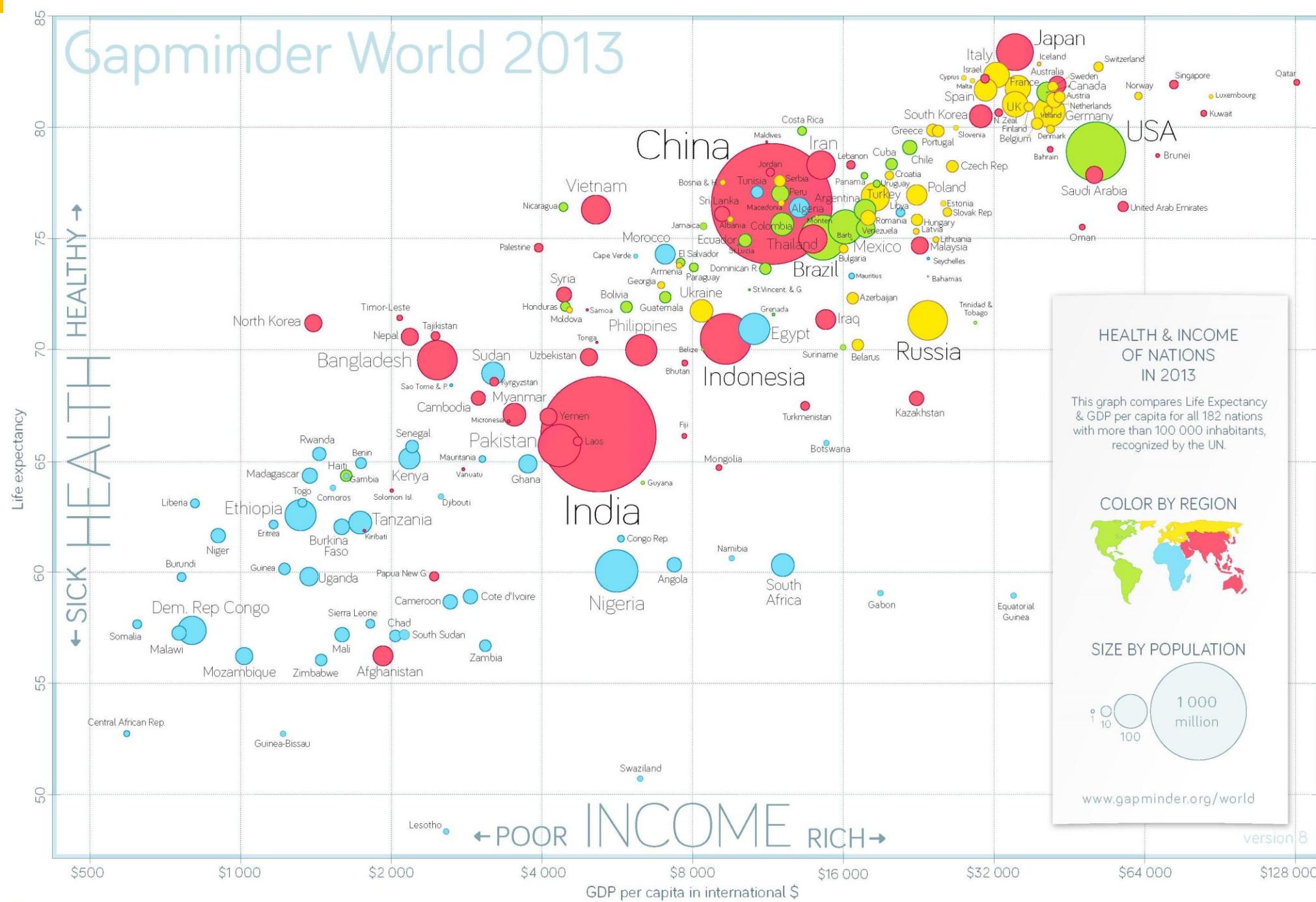
- Accentuate data on scatter plots
- Overlay on maps



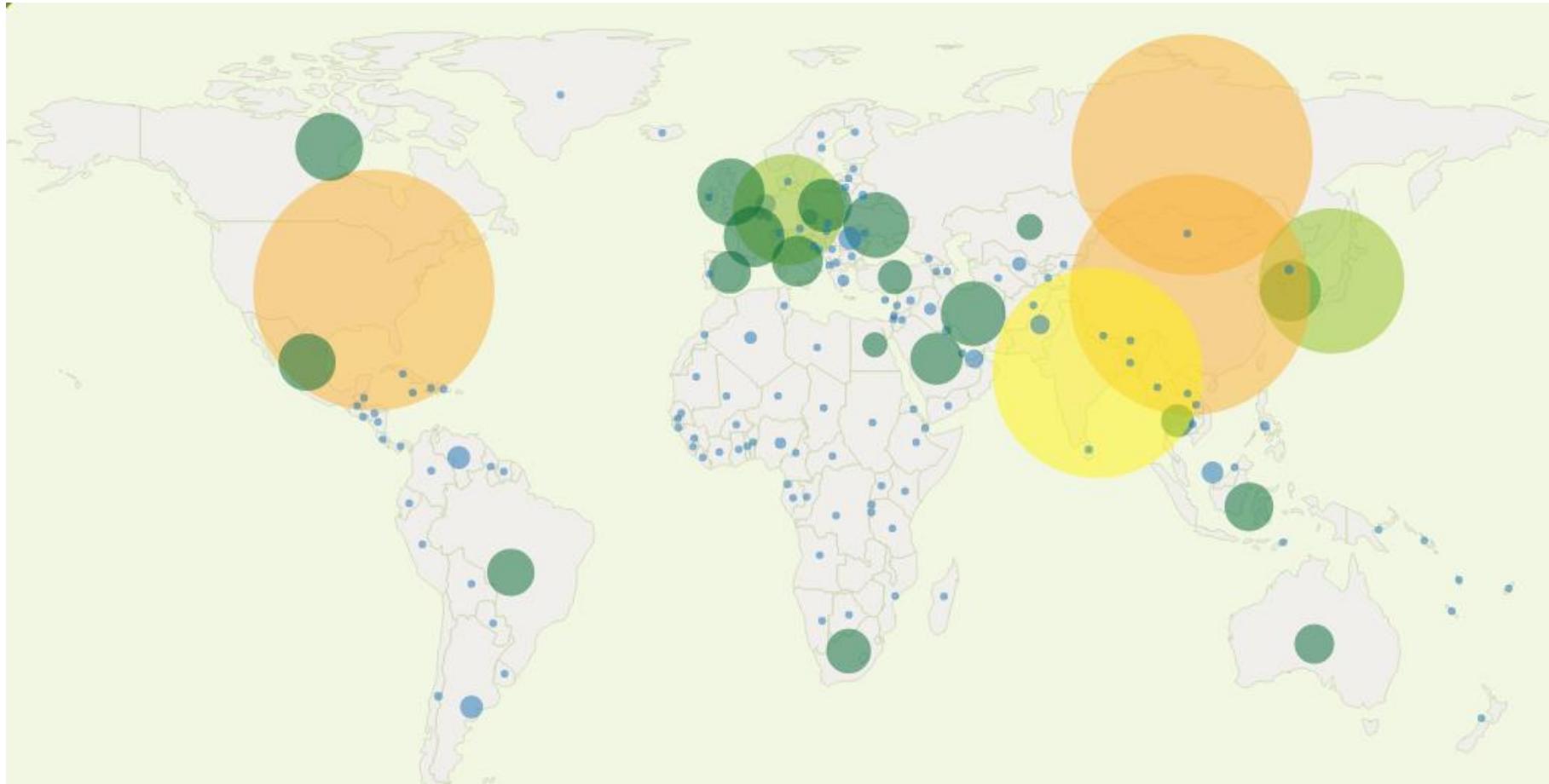


Bubble chart + Scatter plot





Bubble chart + Map



<https://gisgeography.com/climate-change-effects-maps/>

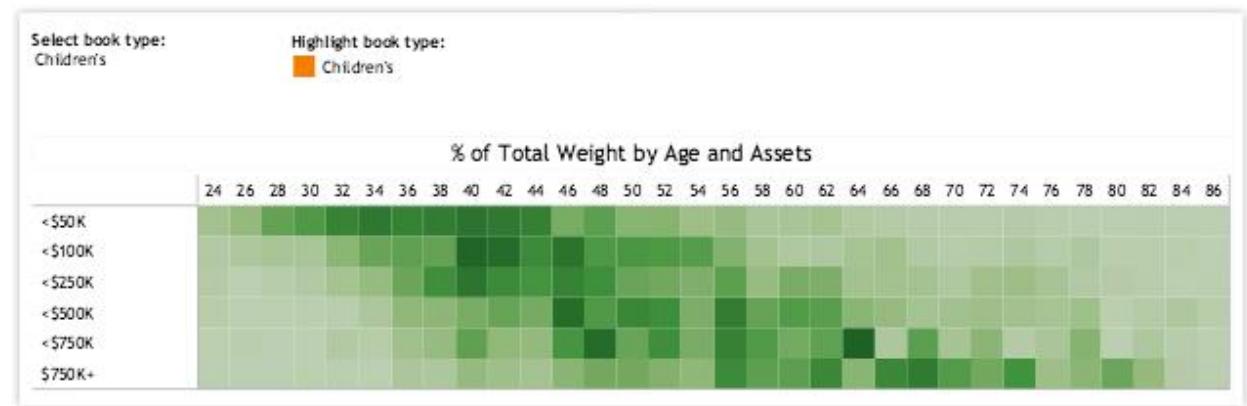
Heatmap

When to use

- Showing the relationship between two factors

Possible extension

- Vary the size of squares
 - Using something other than squares



Crosstab / Highlight table

When to use

- Providing detailed information on heat maps

Possible extension

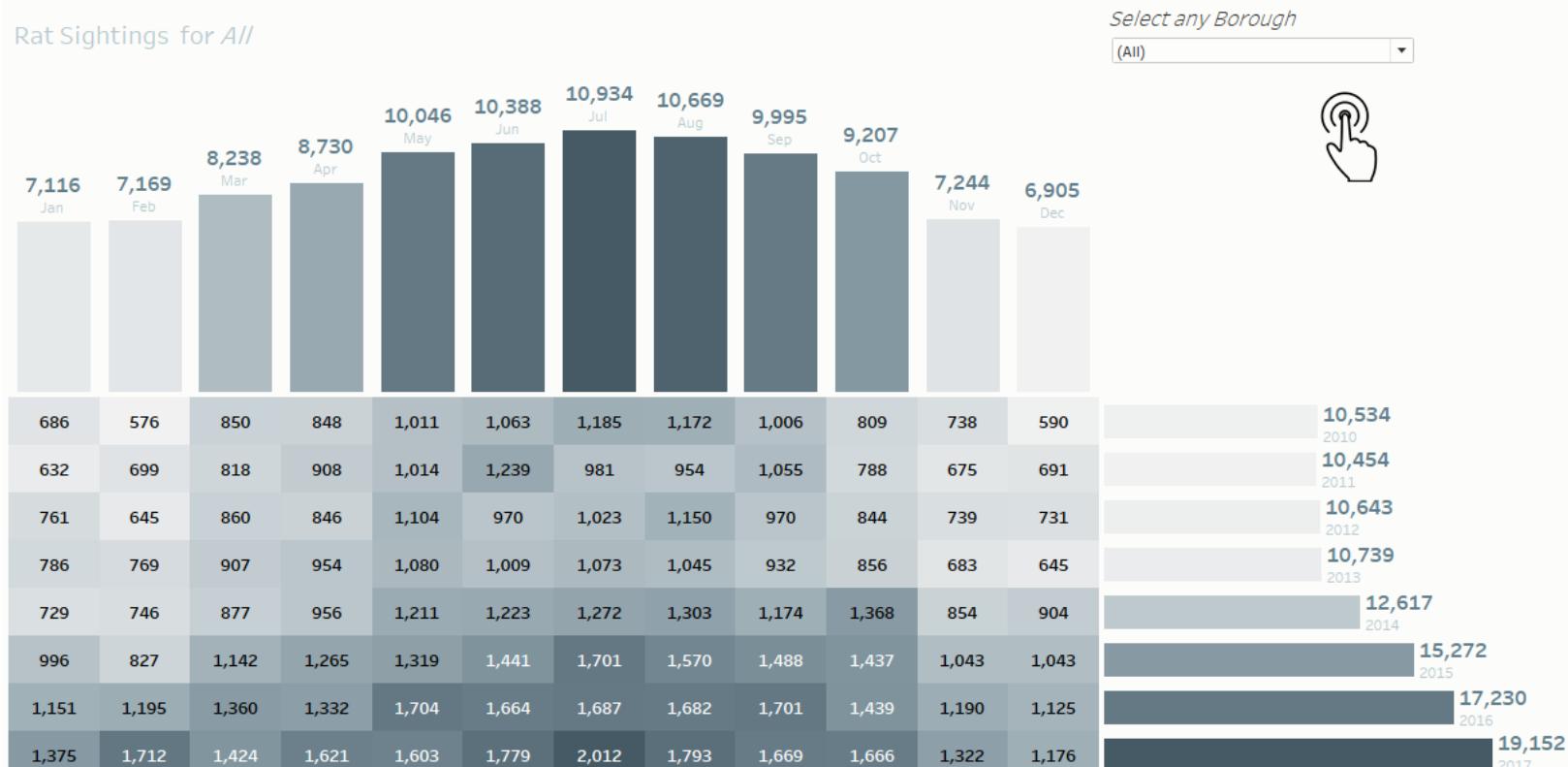
- Combine highlight tables with other chart types

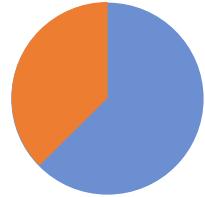
ปีงบประมาณ	January	February	March	April	May	June	July	August	September	October	November	December
2549		1,180,500	305,000	7,500		13,500	10,000	25,000	150,200			
2550	10,500	116,500	1,000		5,500	52,500	1,000	1,500	28,500	2,000	32,500	230,070
2551	1,108,500	411,437	867,134	15,500	101,000	18,000	114,840	53,000	223,500	284,817	640,500	3,946,100
2552	221,500	621,752	59,200	64,000	28,400	6,105,500	86,000	1,087,000	290,500	64,500	95,000	224,899
2553	1,574,100	1,053,856	6,107,843	1,013,515	238,037	1,121,978	87,000	1,635,676	1,514,300	1,546,500	3,145,400	1,078,100
2554	912,500	1,411,624	512,853	465,641	101,500	596,500	1,694,705	203,840	725,833	16,543,300	5,955,120	10,364,730
2555	858,820	924,500	876,300	806,000	724,150	280,550	851,050	404,750	1,134,839	424,133	1,616,928	4,138,347
2556	2,203,699	1,272,799	565,800	905,400	547,651	1,220,000	321,000	982,409	946,698	31,250	2,938,380	1,138,300
2557	737,103	934,225	986,295	2,999,760	958,704	1,411,319	1,291,100	2,613,600	2,169,467	1,267,900	2,050,543	16,888,200
2558	2,645,888	1,907,457	1,054,491	627,722	376,000	1,643,867	3,533,850	16,230,631	11,154,671	952,000	2,970,500	3,027,865
2559	2,797,940	2,864,000	14,684,125	2,900,500	483,635	2,590,946	2,442,828	5,484,750	2,236,651	1,334,602	1,635,100	12,888,969
2560	2,751,941	1,558,290	4,560,788	478,278	937,405	3,218,400	1,336,045	3,019,639	2,106,749	3,051,750	3,466,650	4,965,579
2561	418,440	1,634,557	6,898,784	3,238,000	458,700	1,862,243	900,900	2,300,950	3,926,657	2,511,812	2,619,869	5,045,770
2562	1,163,000	1,736,500	2,510,200	1,718,308	1,433,568	578,959				5,264,100	6,816,792	5,958,382

Crosstab + Bar chart

Rat Sightings in New York City

Since 2010, there have been an increasing number of reported rat sightings across the five NYC Boroughs, peaking each year in the summer months.





3. Composition

Pie chart

Tree map

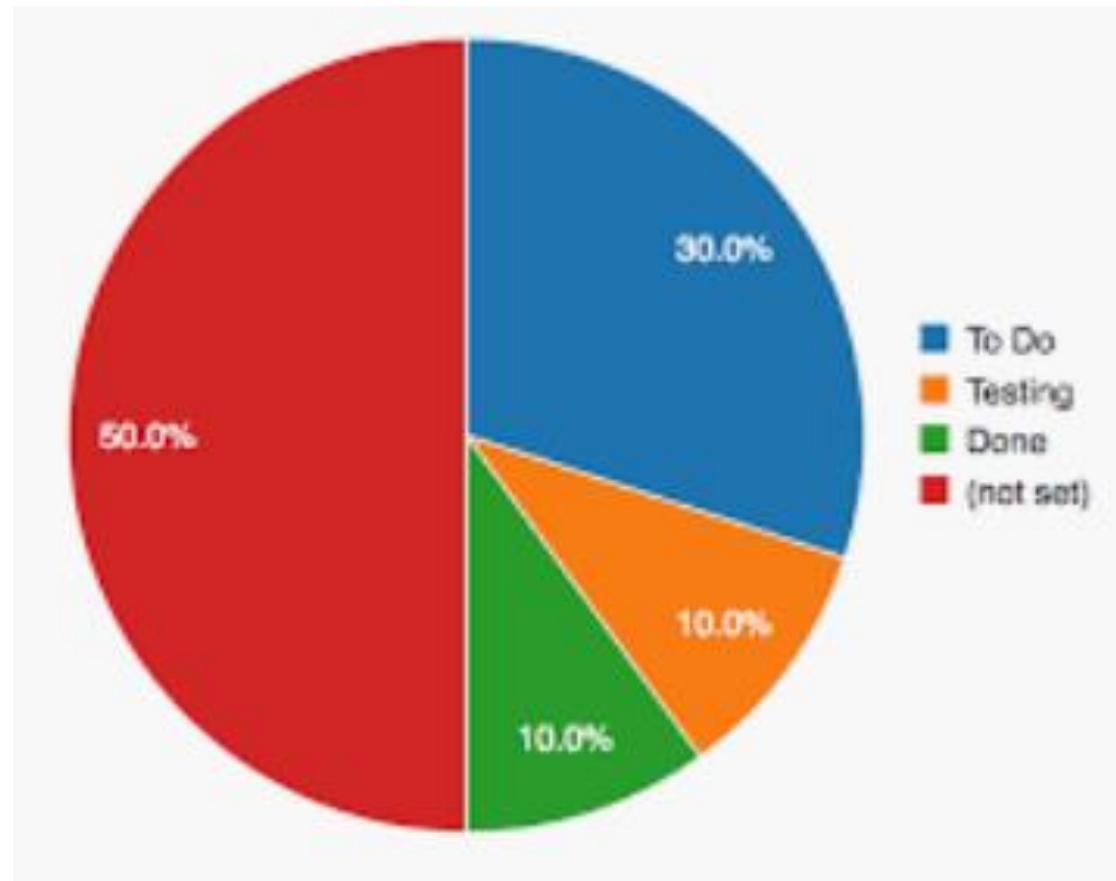
Pie chart

When to use

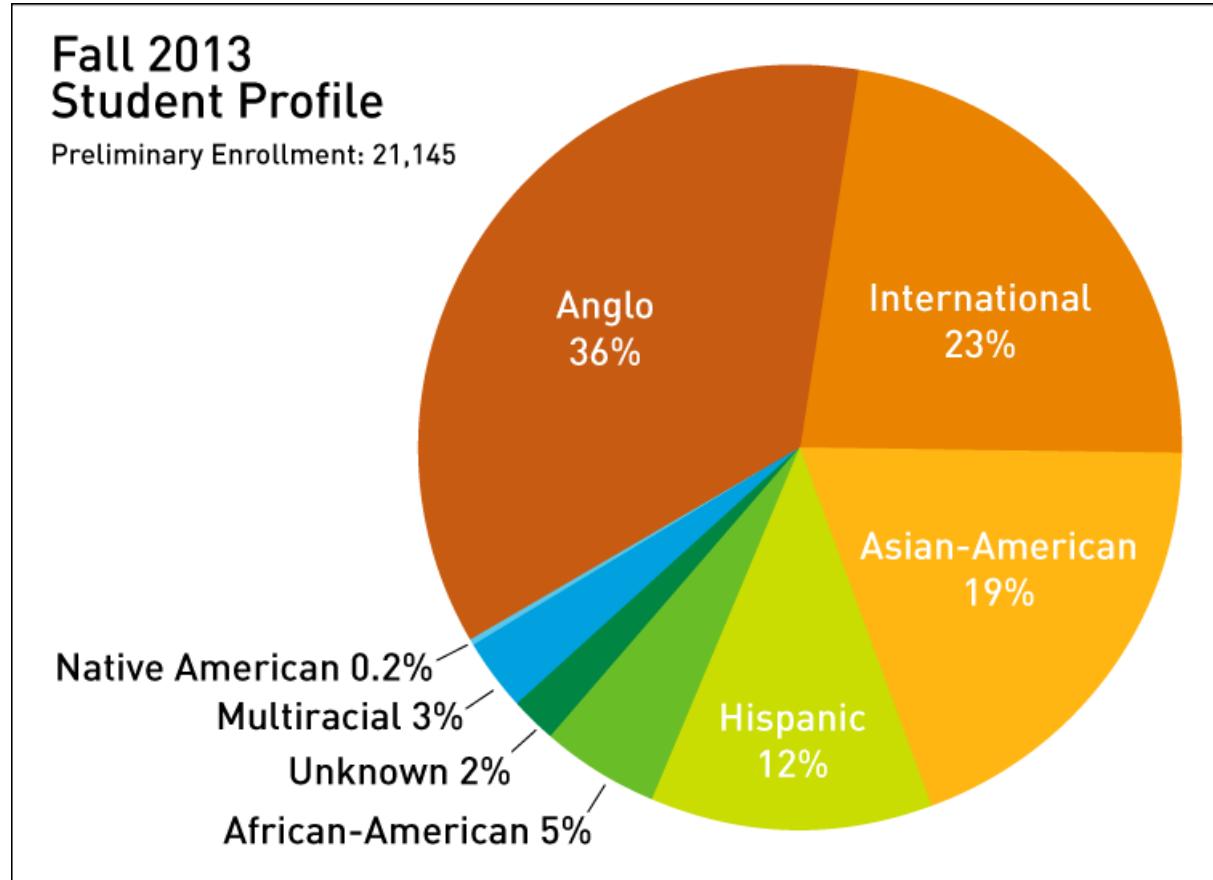
- Showing proportions

Possible extension

- Limit pie wedges to six
- Overlay pies on maps



Pie chart



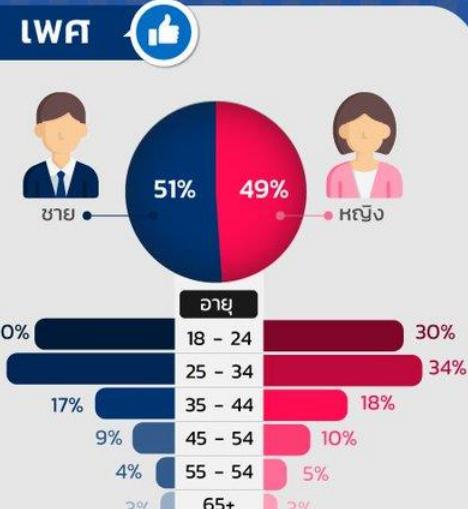
- UT Dallas' diversity breakdown
- Total of 100.2 !

สถิติของประชากร Facebook ในไทย



ประชากร Facebook ในไทยก้าวหน้าไปประมาณ 50 ล้านบัญชี

เพศ



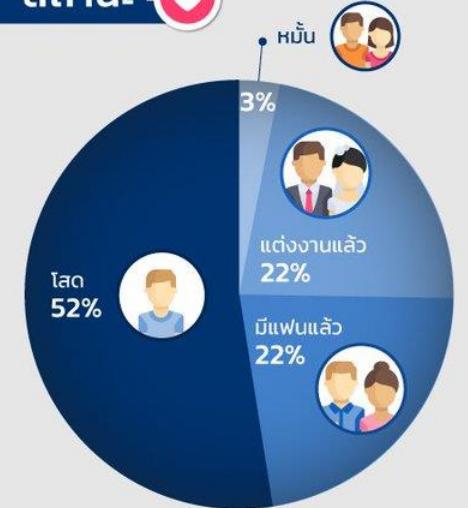
"กลุ่มคนวัยทำงาน" เล่น Facebook มากที่สุด

การศึกษา



ส่วนใหญ่จบการศึกษา "ปริญญาตรี"

สถานะ



ครึ่งหนึ่งของคนเล่น Facebook มีสถานะ "โสด"

อาชีพ

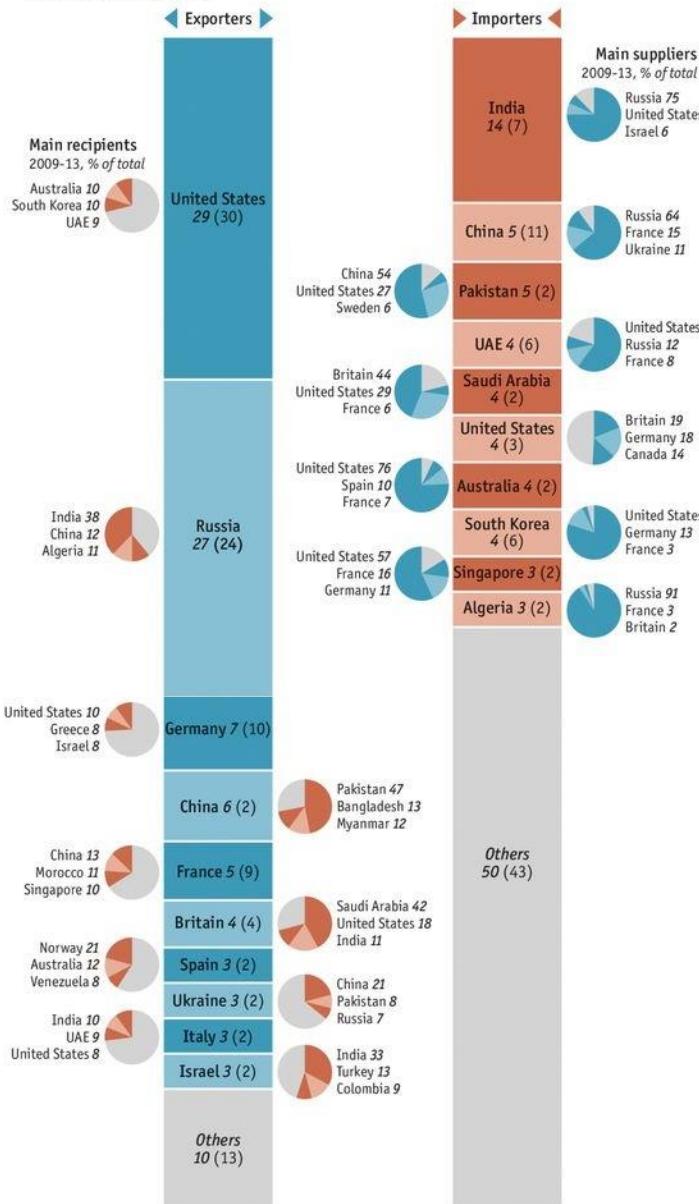
10 อันดับแรก



อาชีพยอดนิยม คือ "ผู้บริหาร"

International arms sales

2009-13 (2004-08), % of total

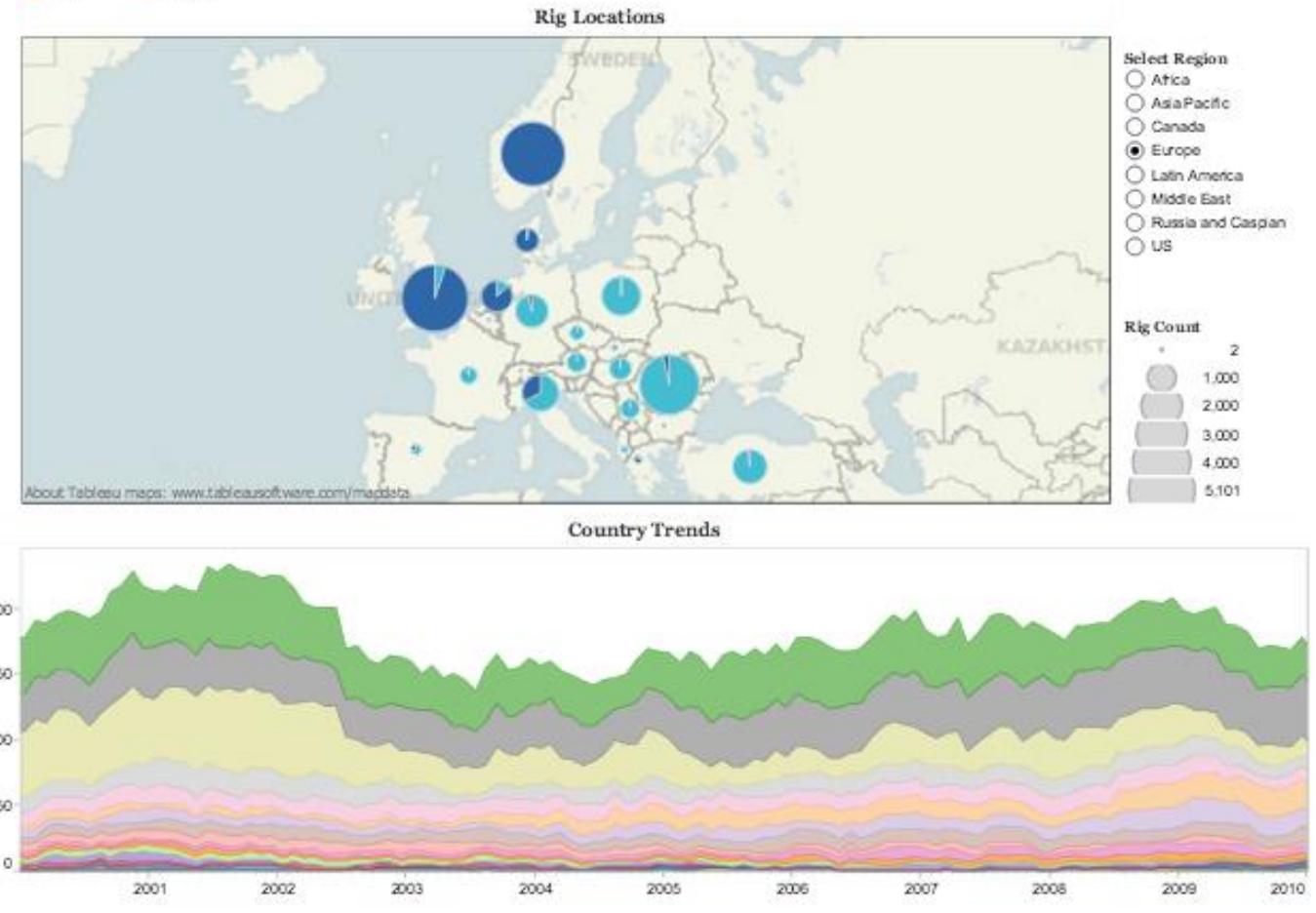


Source: SIPRI

Economist.com/graphicdetail

Worldwide Oil Rigs

Land Offshore



Treemap

When to use

- Showing hierarchical data as a proportion of a whole

Possible extension

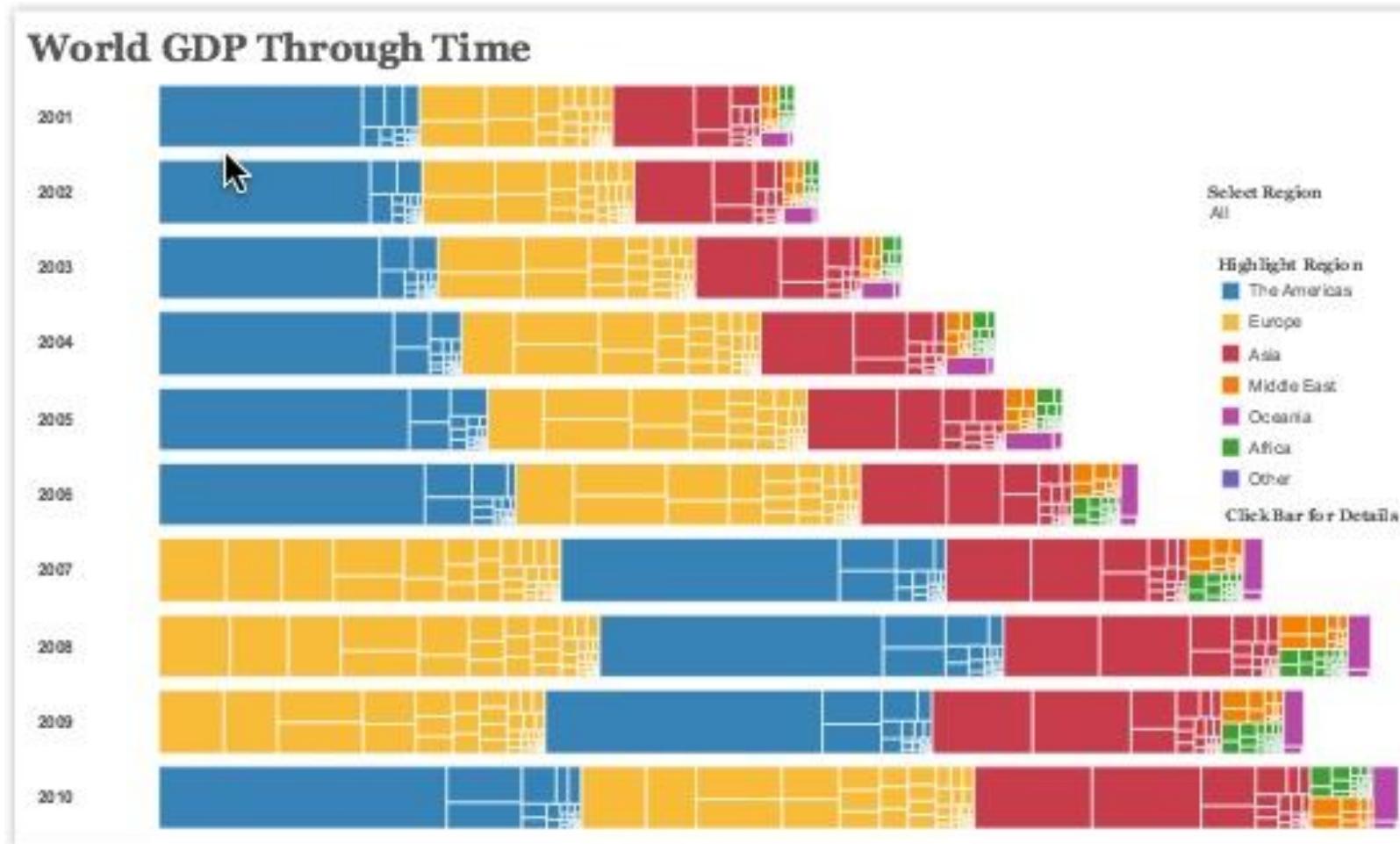
- Colouring the rectangles by a category
- Combining tree maps with bar charts

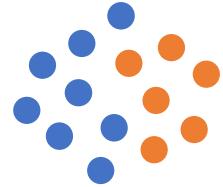


Treemap + Coloring



Treemap + Bar chart





4. Distribution

Histogram chart

Box plot

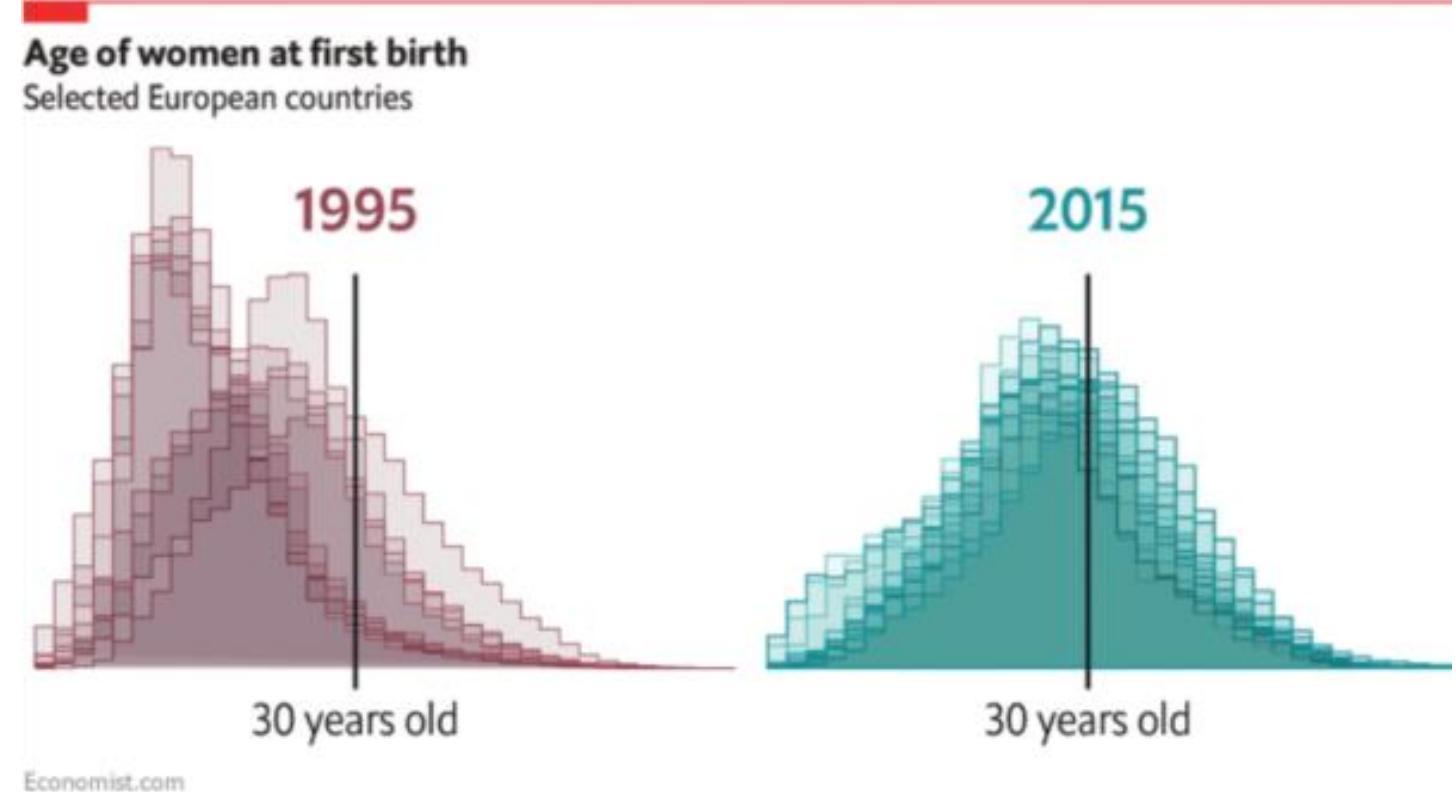
Histogram chart

When to use

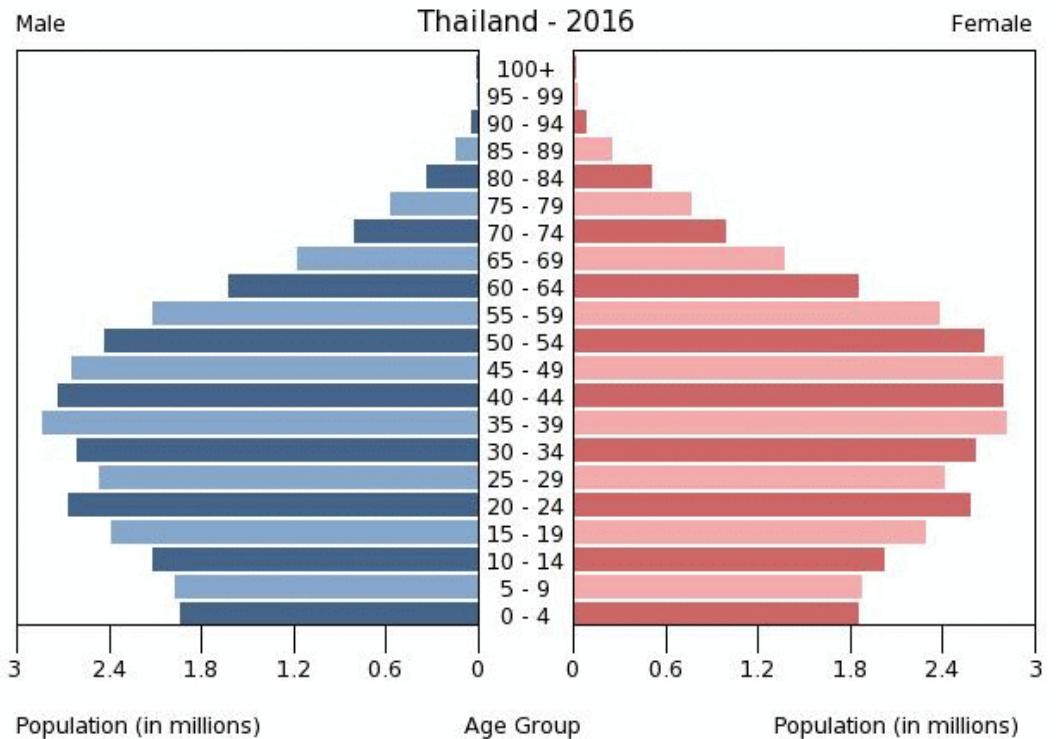
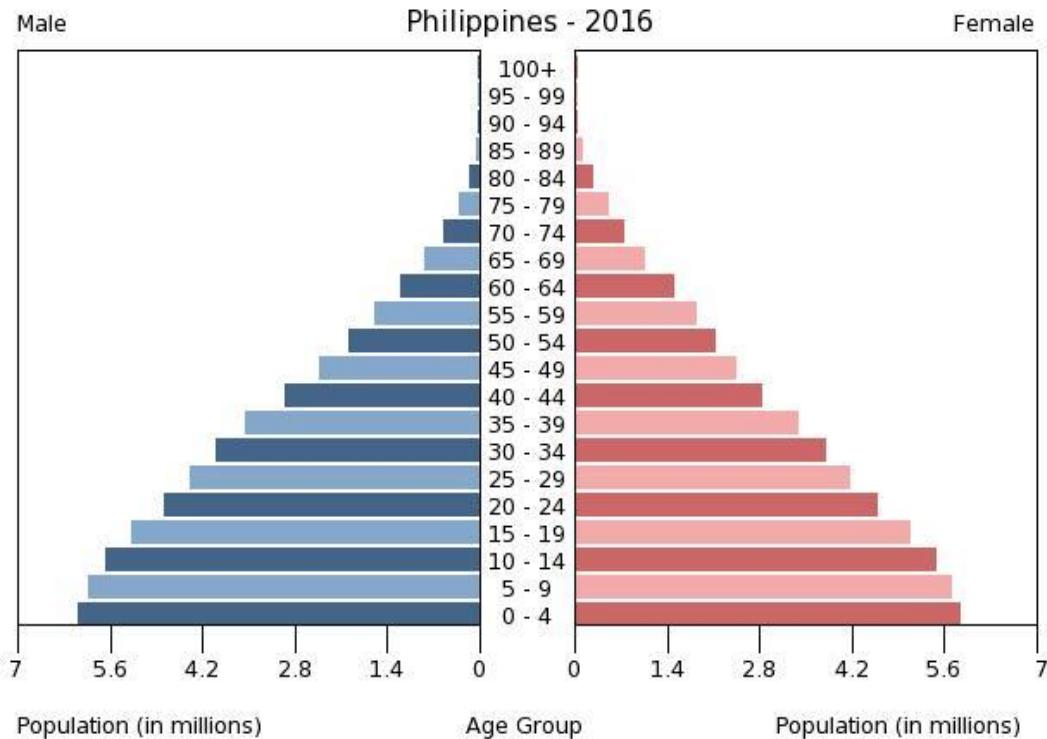
- Understanding the distribution of your data

Possible extension

- Test different groupings of data
 - Add a filter



Histogram chart



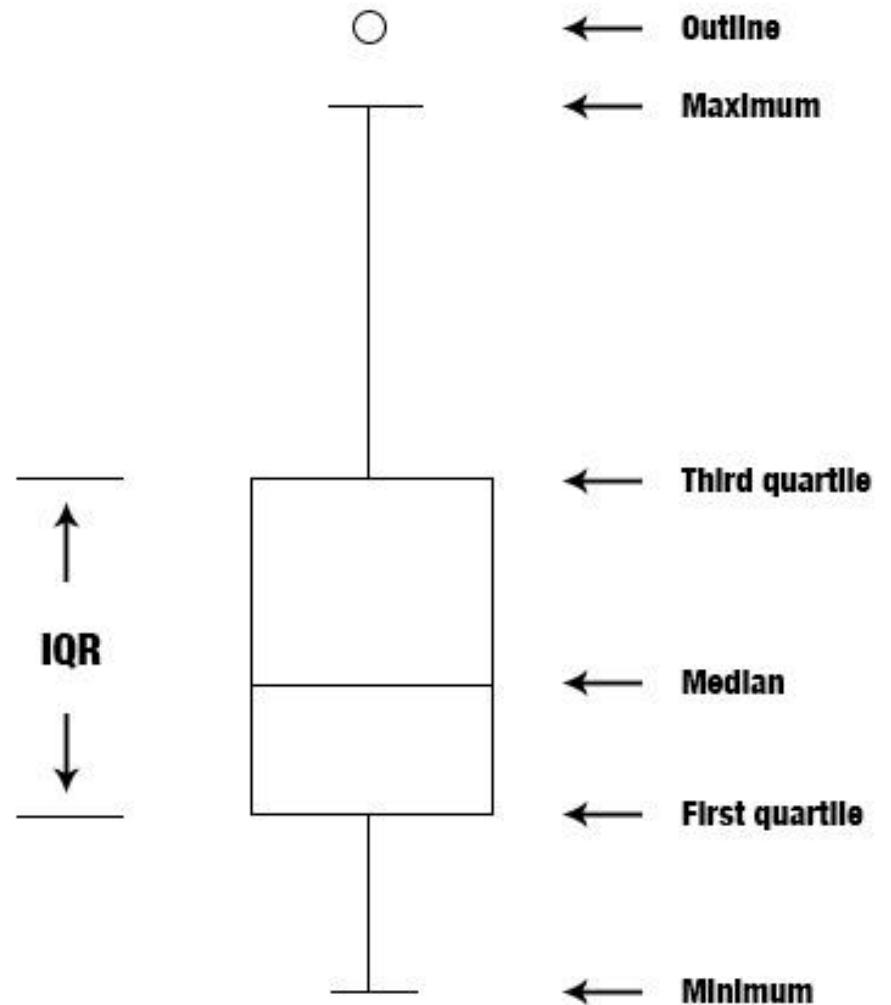
Box plot

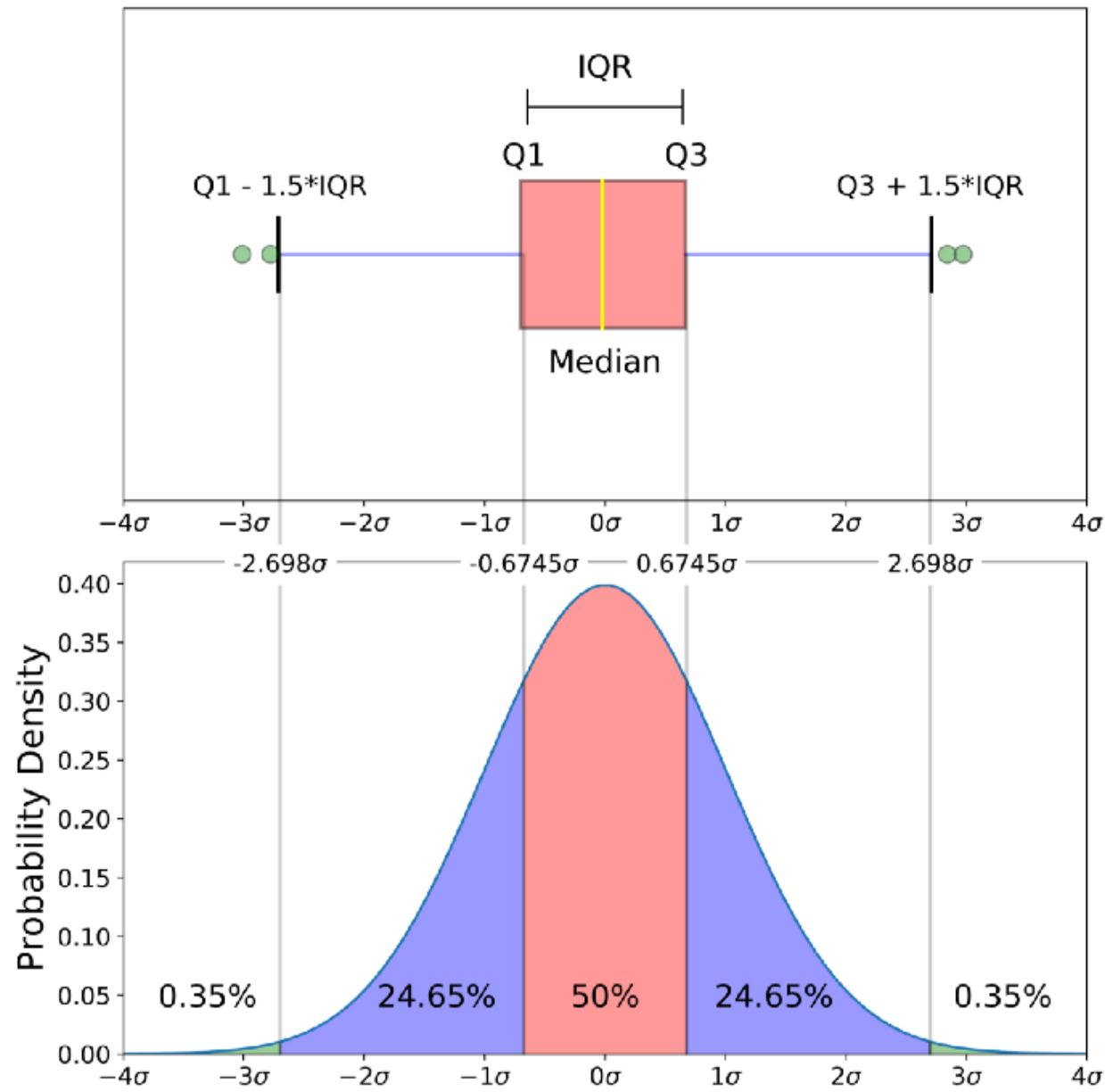
When to use

- Showing the distribution of a set of data

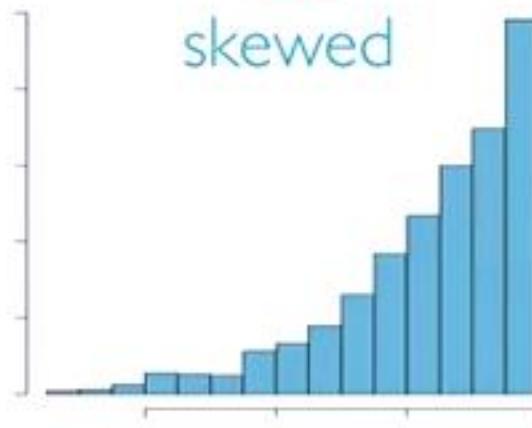
Possible extension

- Hiding the points within the box
- Comparing boxplots across categorical dimensions

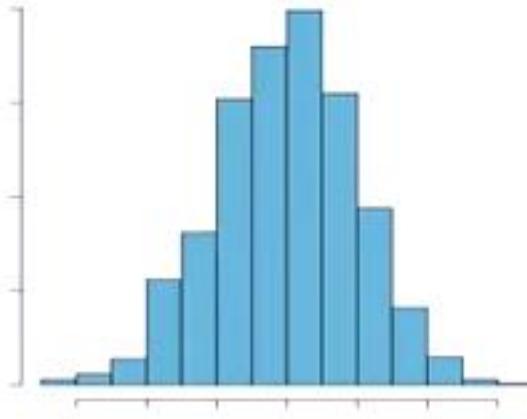




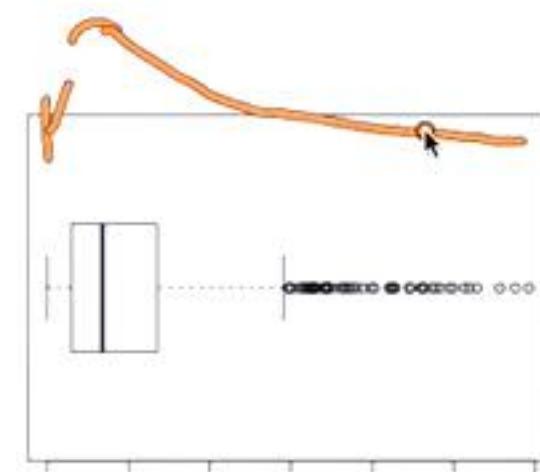
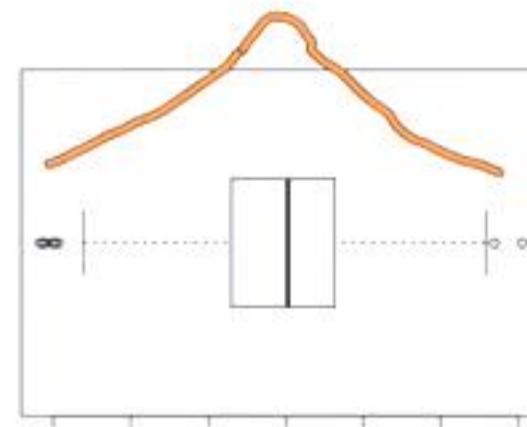
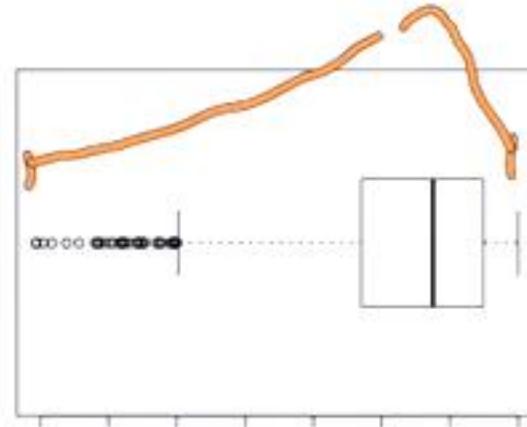
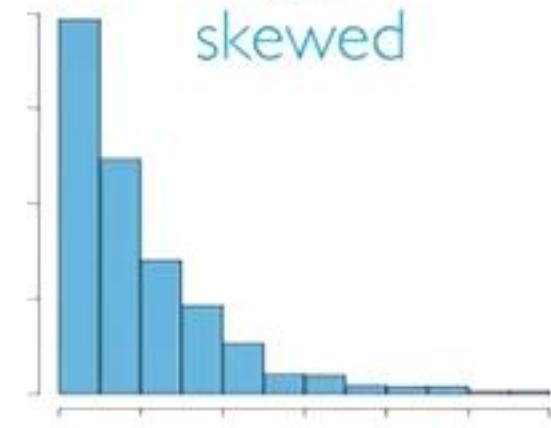
left
skewed



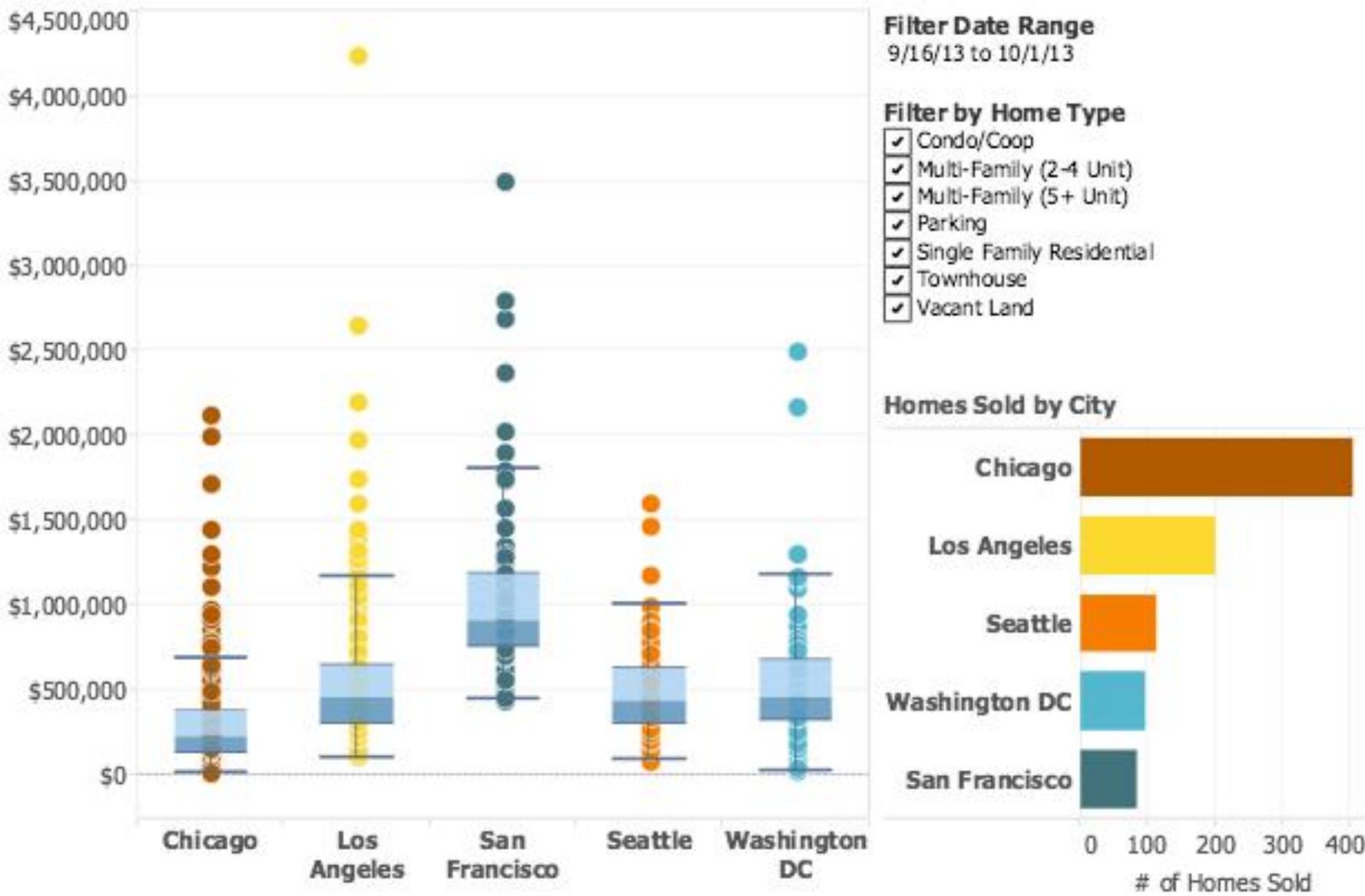
symmetric



right
skewed



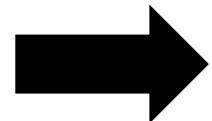
Two Weeks of Home Sales



What do you want to tell from your data?



Story



1.Comparison

Bar chart, Line chart
Bullet chart



2.Relationship

Scatter plot, Map
Bubble chart, Heat map
Crosstab / Highlight table



3.Composition

Pie chart
Tree map

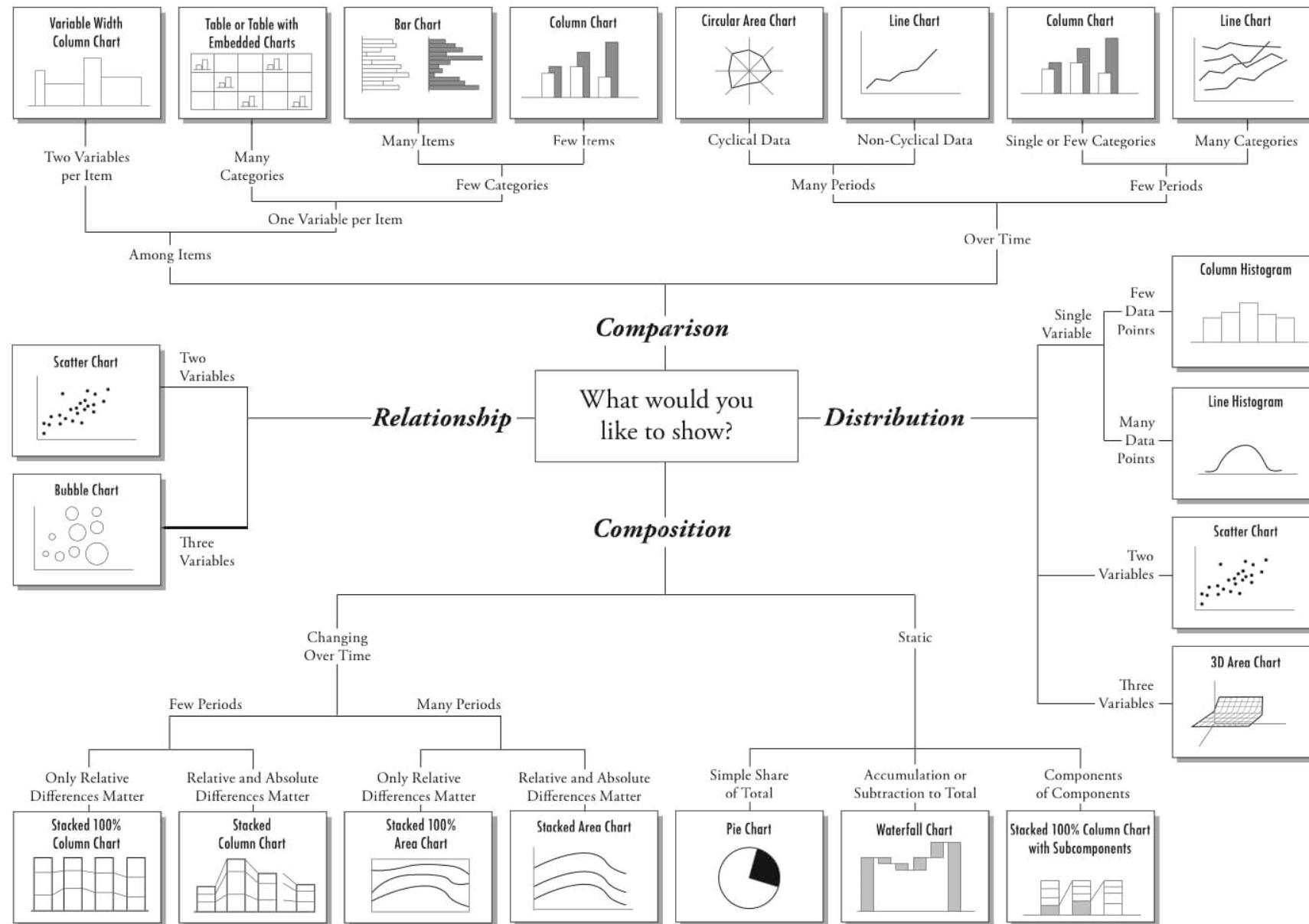


4.Distribution

Histogram chart
Box plot

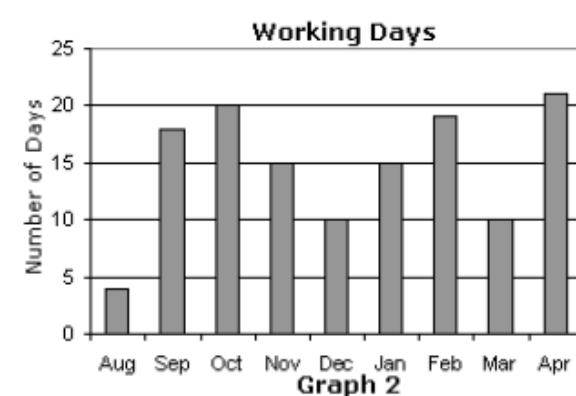
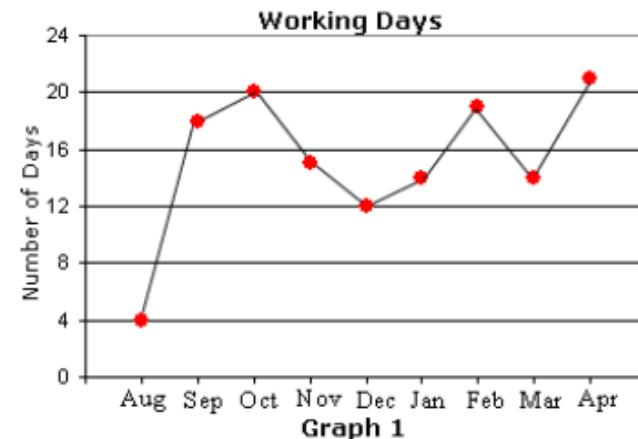
Chart Suggestions—A Thought-Starter

www.ExtremePresentation.com
 © 2009 A. Abela — a.abela@gmail.com



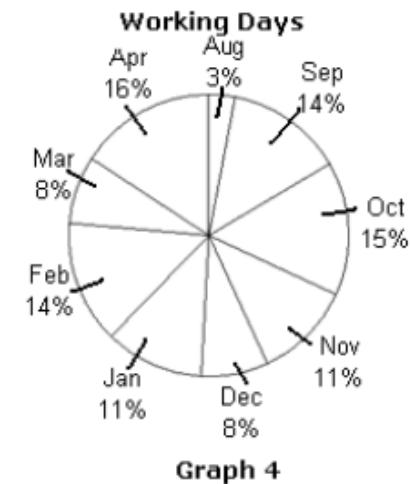
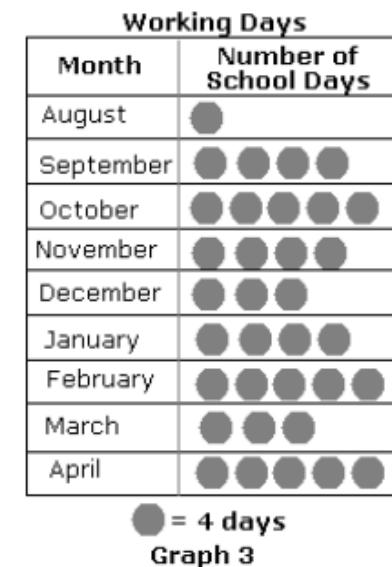
Month	Aug	Sep	Oct	Nov	Dec	Jan	Feb	Mar	Apr
Number of Days	4	18	20	15	12	14	19	14	21

Relationship >

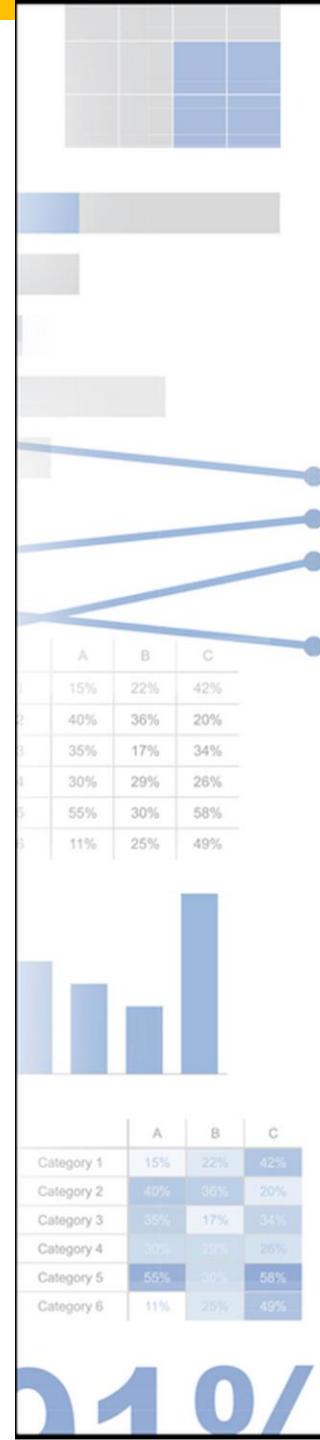


< Comparison

Distribution >



< Composition



cole nussbaumer knaflic

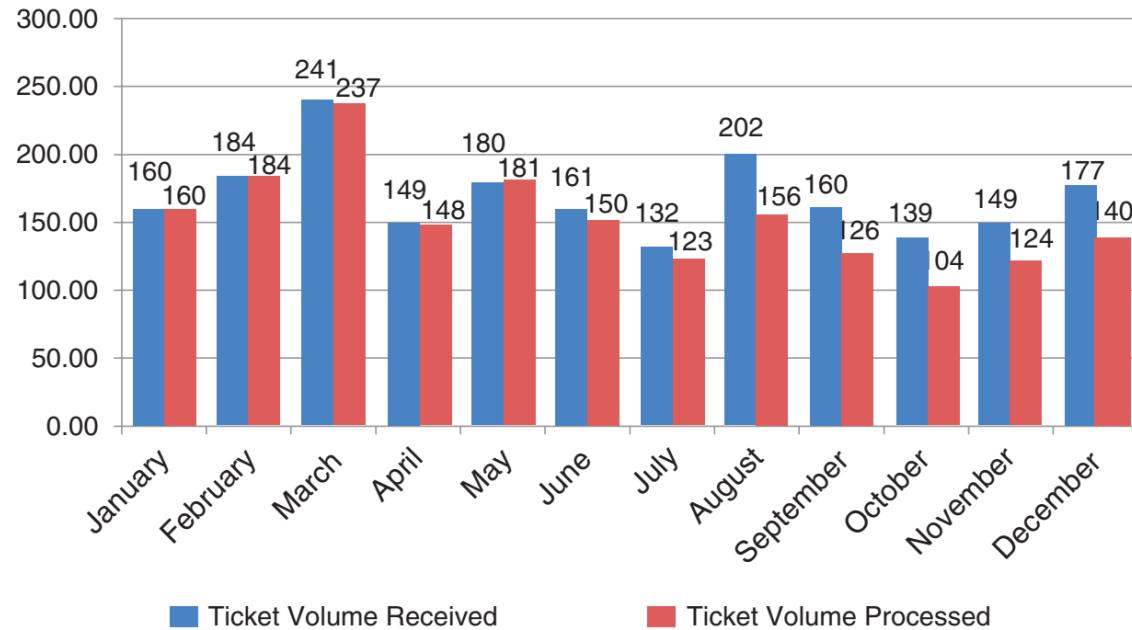
storytelling with data

a data
visualization
guide for
business
professionals

010/

WILEY

Ticket Trend



Ticket volume over time

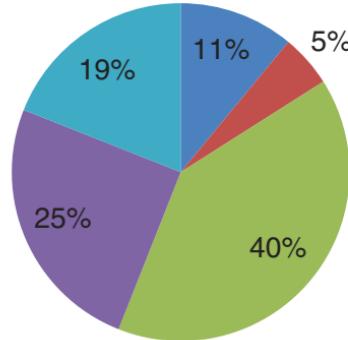


Data source: XYZ Dashboard, as of 12/31/2014 | A detailed analysis on tickets processed per person and time to resolve issues was undertaken to inform this request and can be provided if needed.

Survey Results

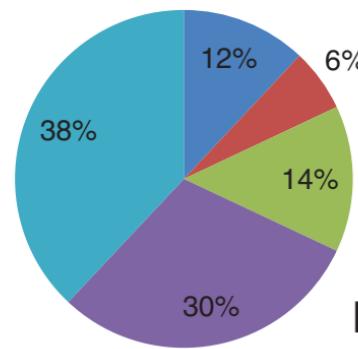
PRE: How do you feel about doing science?

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



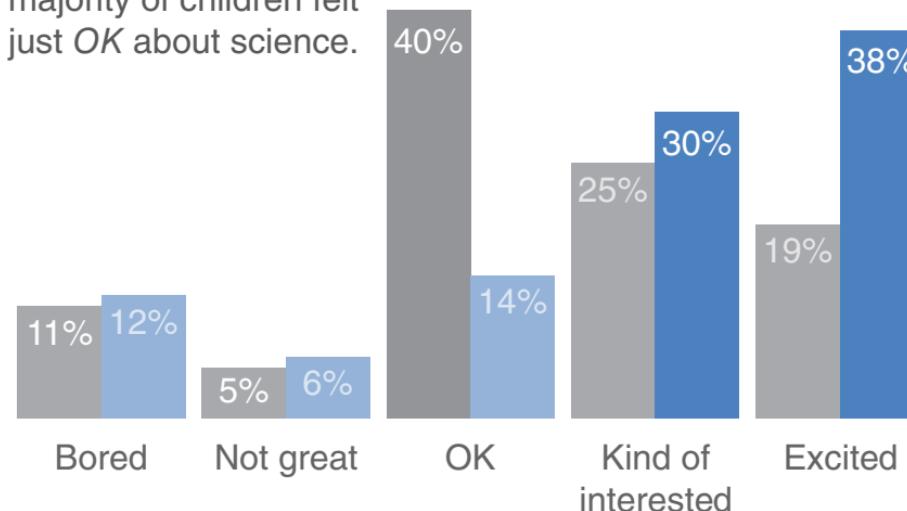
POST: How do you feel about doing science?

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



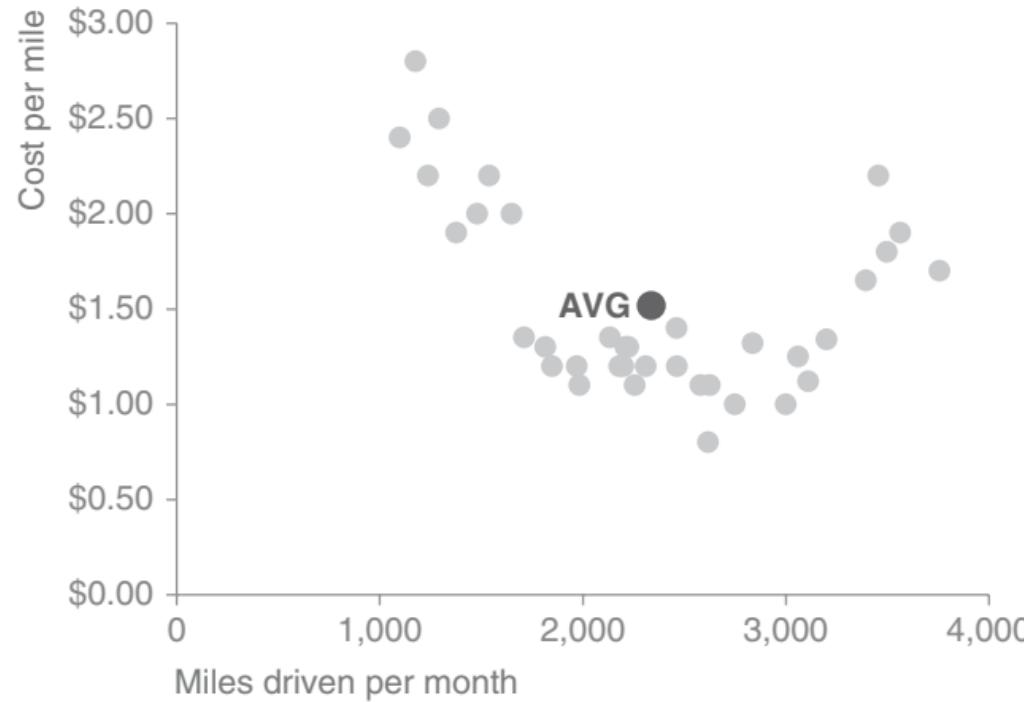
How do you feel about science?

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

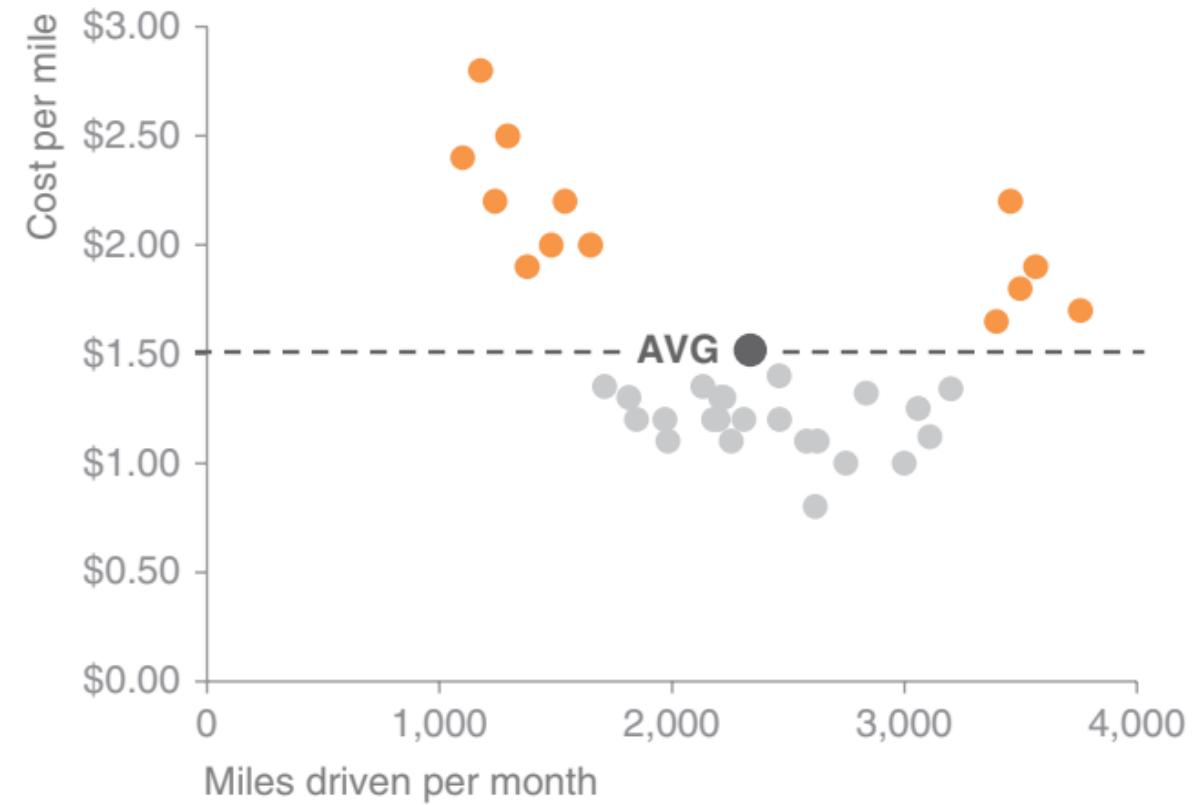


AFTER program, more children were *Kind of interested* & *Excited* about science.

Cost per mile by miles driven

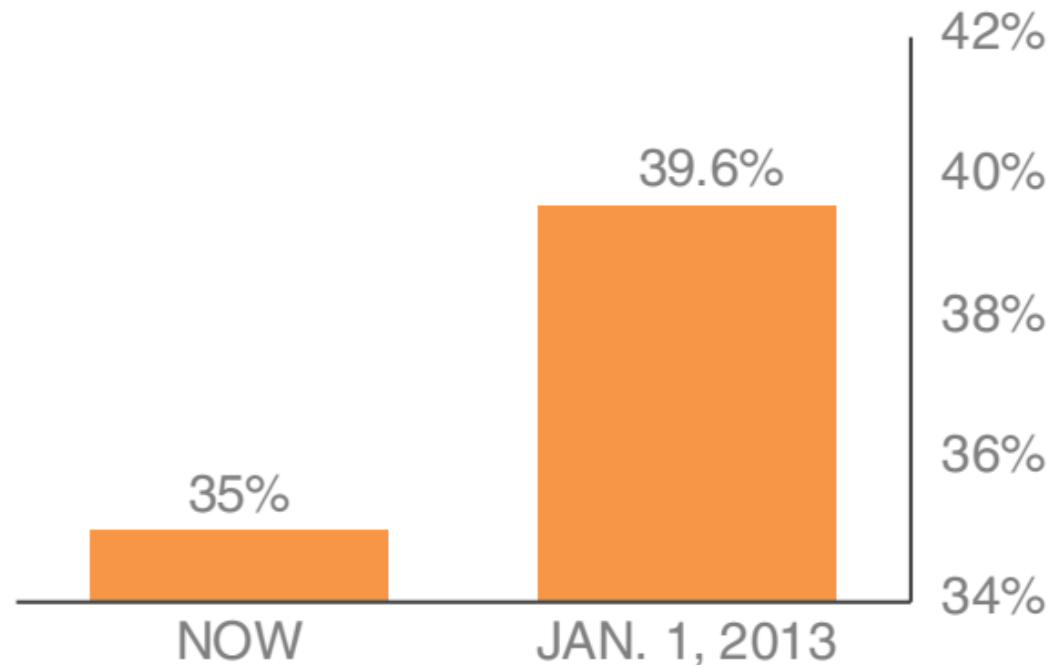


Cost per mile by miles driven



Non-zero baseline: as originally graphed

IF BUSH TAX CUTS EXPIRE
TOP TAX RATE



Zero baseline: as it should be graphed

IF BUSH TAX CUTS EXPIRE
TOP TAX RATE

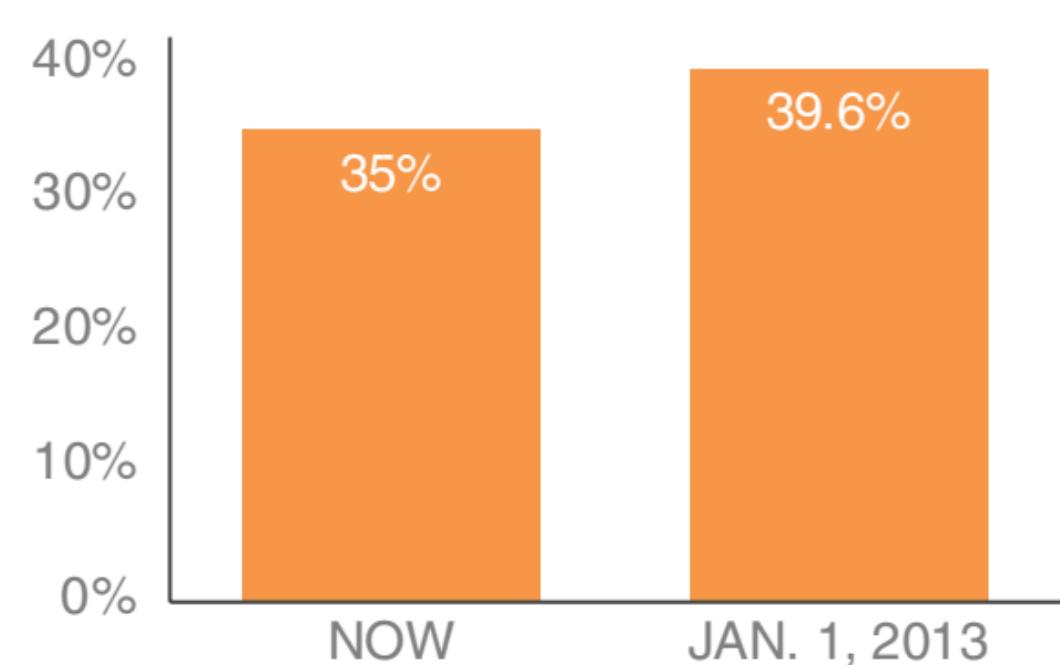
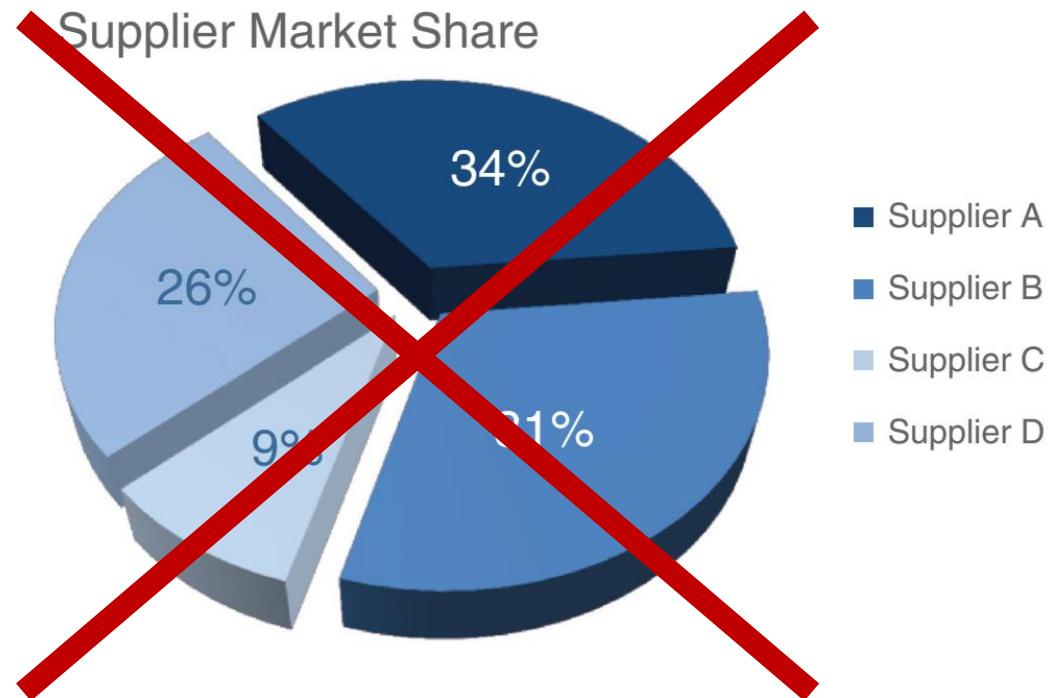


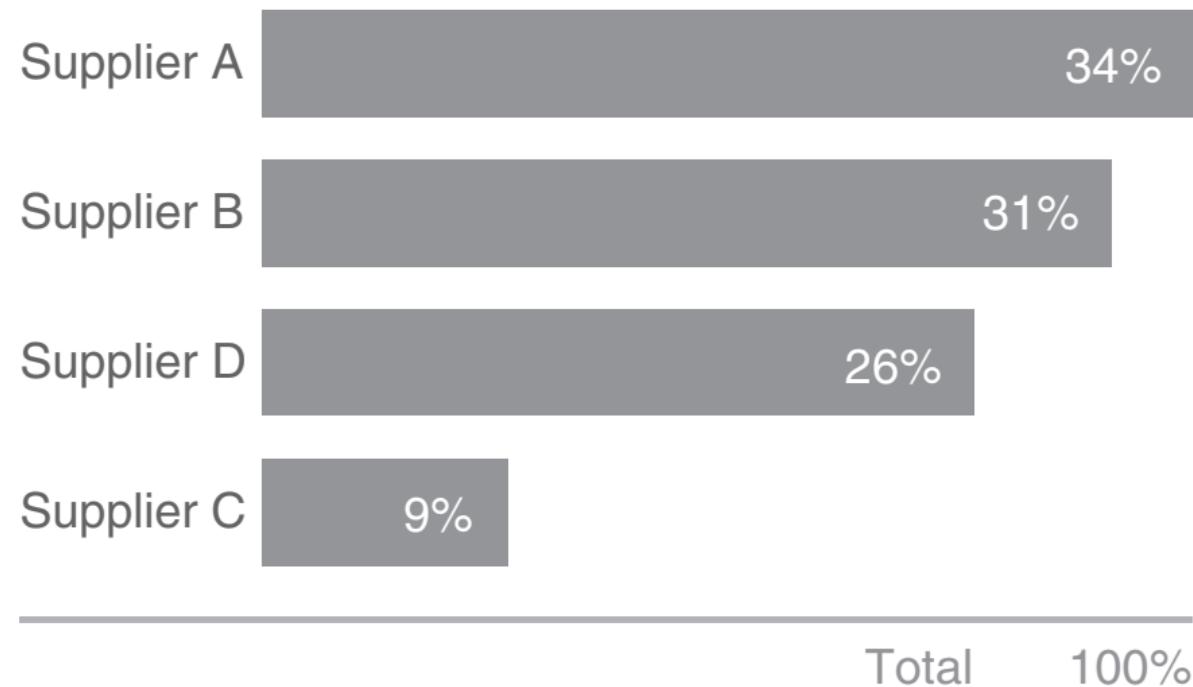
FIGURE 2.13 Bar charts must have a zero baseline

Supplier Market Share

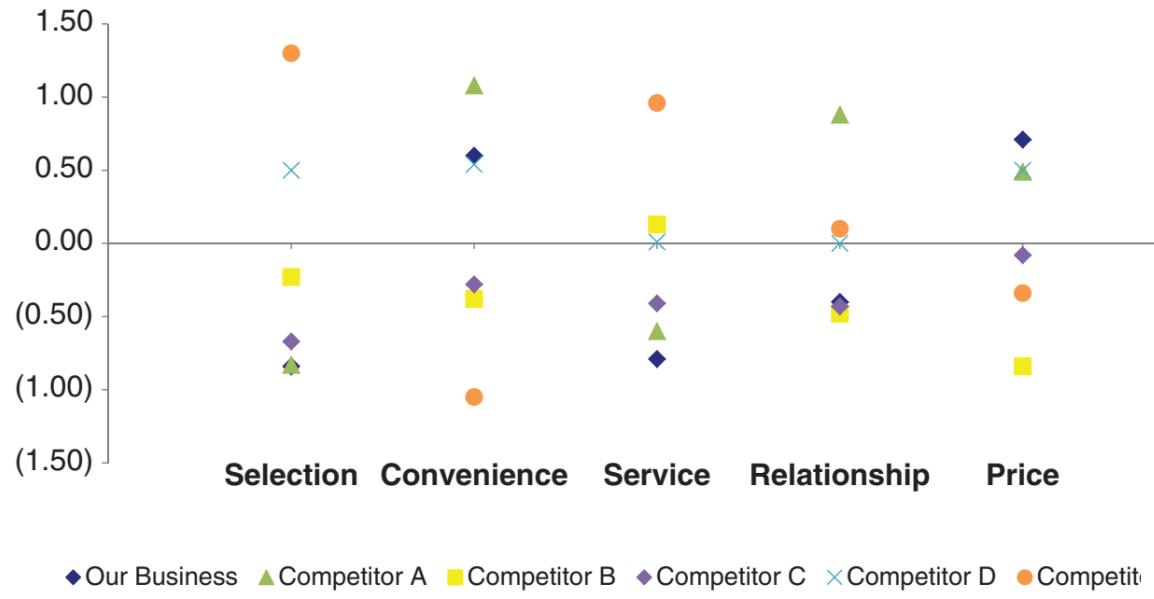


- Supplier A
- Supplier B
- Supplier C
- Supplier D

Supplier Market Share

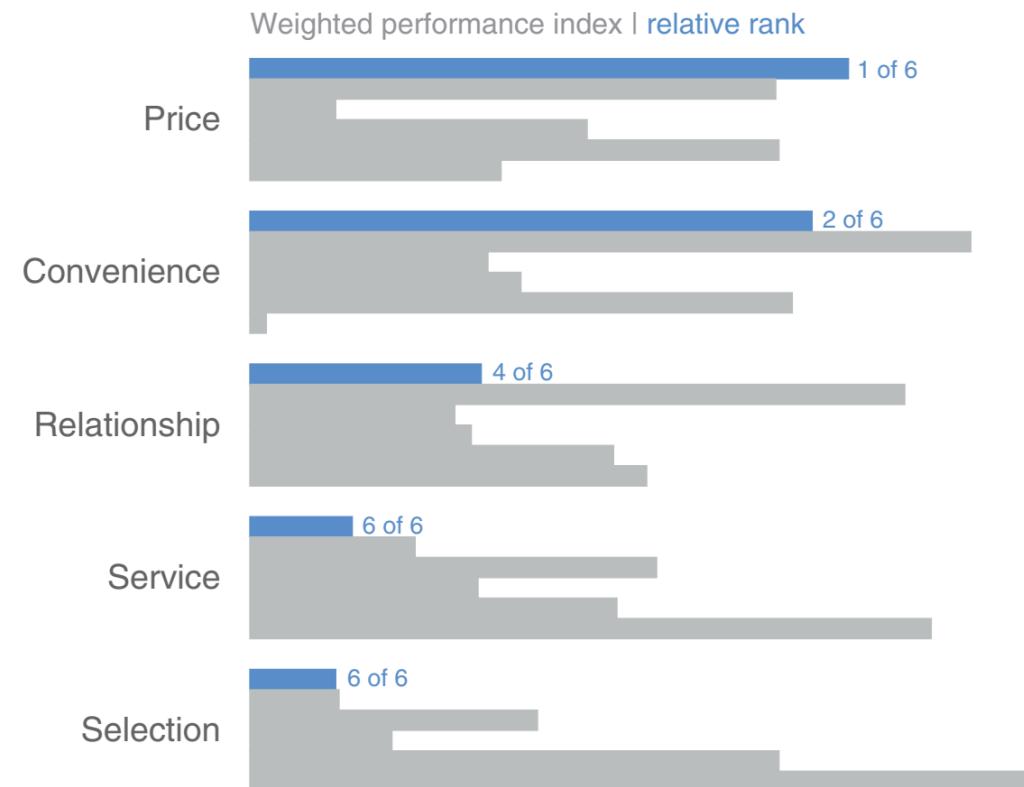


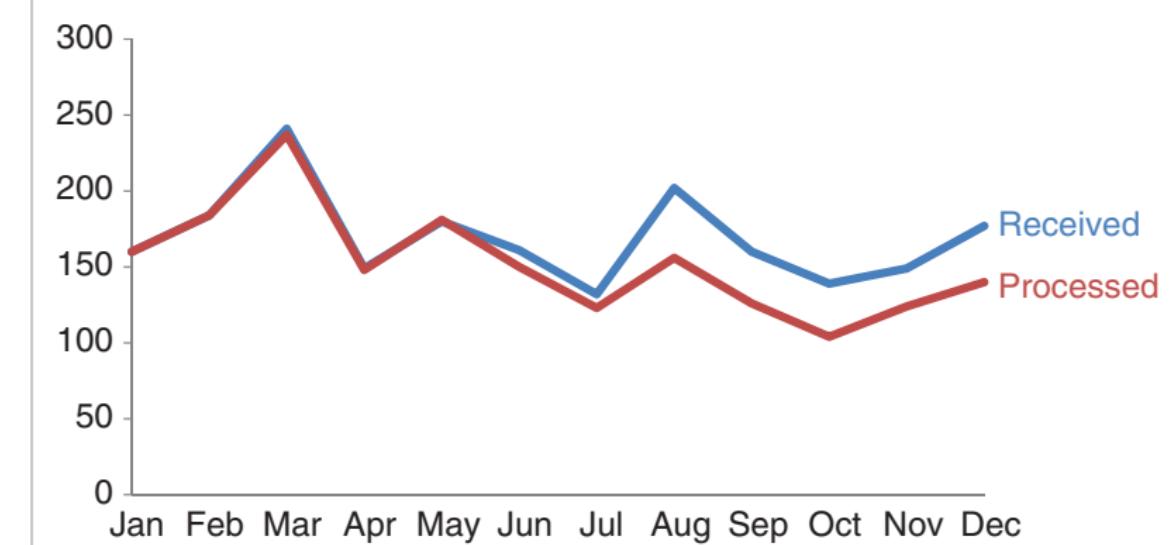
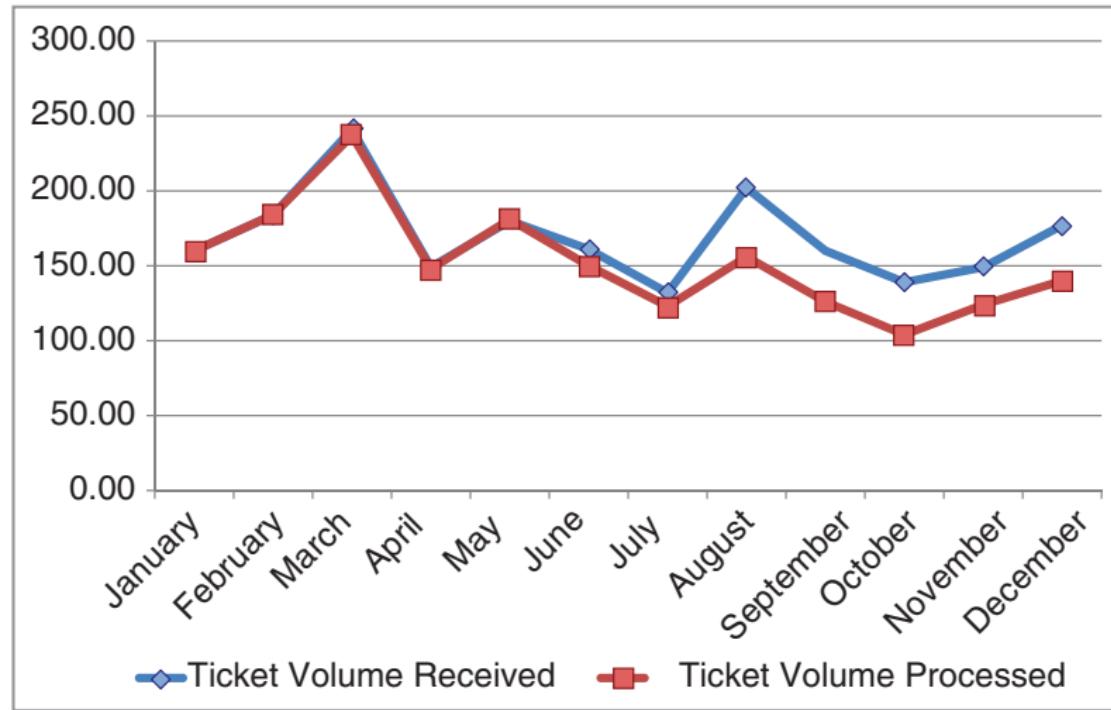
Weighted Performance Index



Performance overview

Weighted performance index | relative rank





Country Level Sales Rank Top 5 Drugs

Rainbow distribution in color indicates sales rank in given country from #1 (red) to #10 or higher (dark purple)

Country	A	B	C	D	E
AUS	1	2	3	6	7
BRA	1	3	4	5	6
CAN	2	3	6	12	8
CHI	1	2	8	4	7
FRA	3	2	4	8	10
GER	3	1	6	5	4
IND	4	1	8	10	5
ITA	2	4	10	9	8
MEX	1	5	4	6	3
RUS	4	3	7	9	12
SPA	2	3	4	5	11
TUR	7	2	3	4	8
UK	1	2	3	6	7
US	1	2	4	3	5

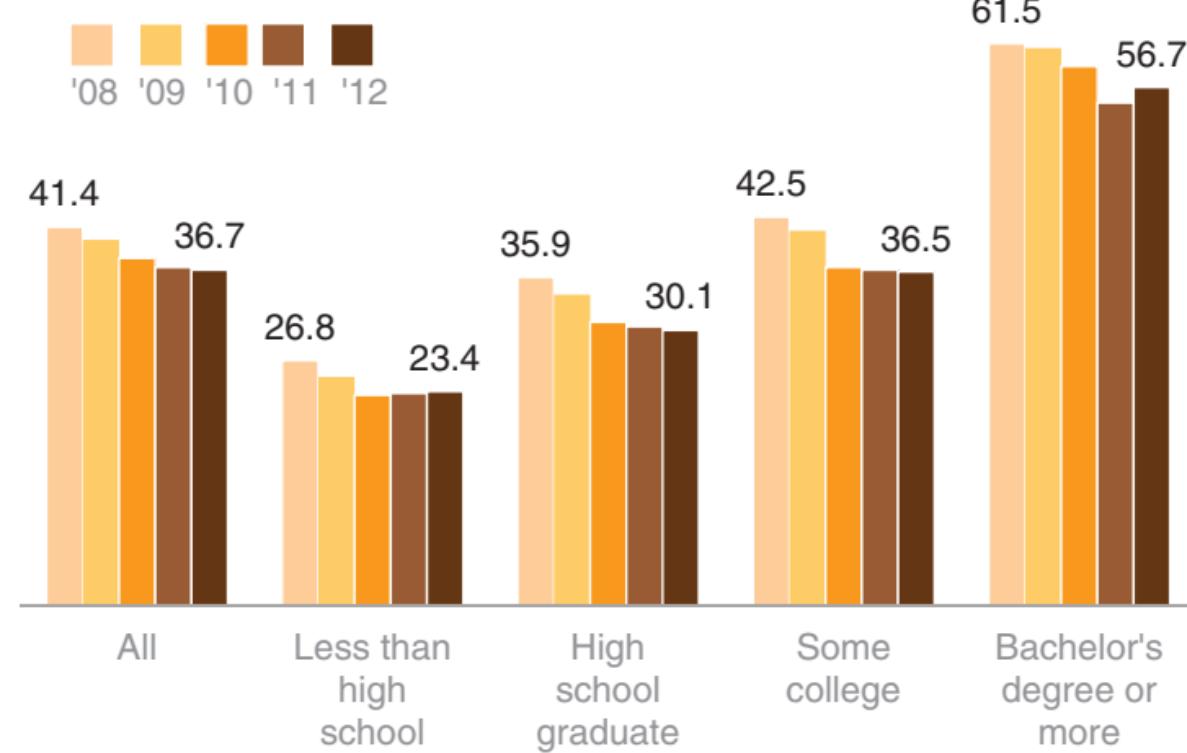
Top 5 drugs: country-level sales rank

RANK	1	2	3	4	5+
COUNTRY DRUG	A	B	C	D	E

Australia	1	2	3	6	7
Brazil	1	3	4	5	6
Canada	2	3	6	12	8
China	1	2	8	4	7
France	3	2	4	8	10
Germany	3	1	6	5	4
India	4	1	8	10	5
Italy	2	4	10	9	8
Mexico	1	5	4	6	3
Russia	4	3	7	9	12
Spain	2	3	4	5	11
Turkey	7	2	3	4	8
United Kingdom	1	2	3	6	7
United States	1	2	4	3	5

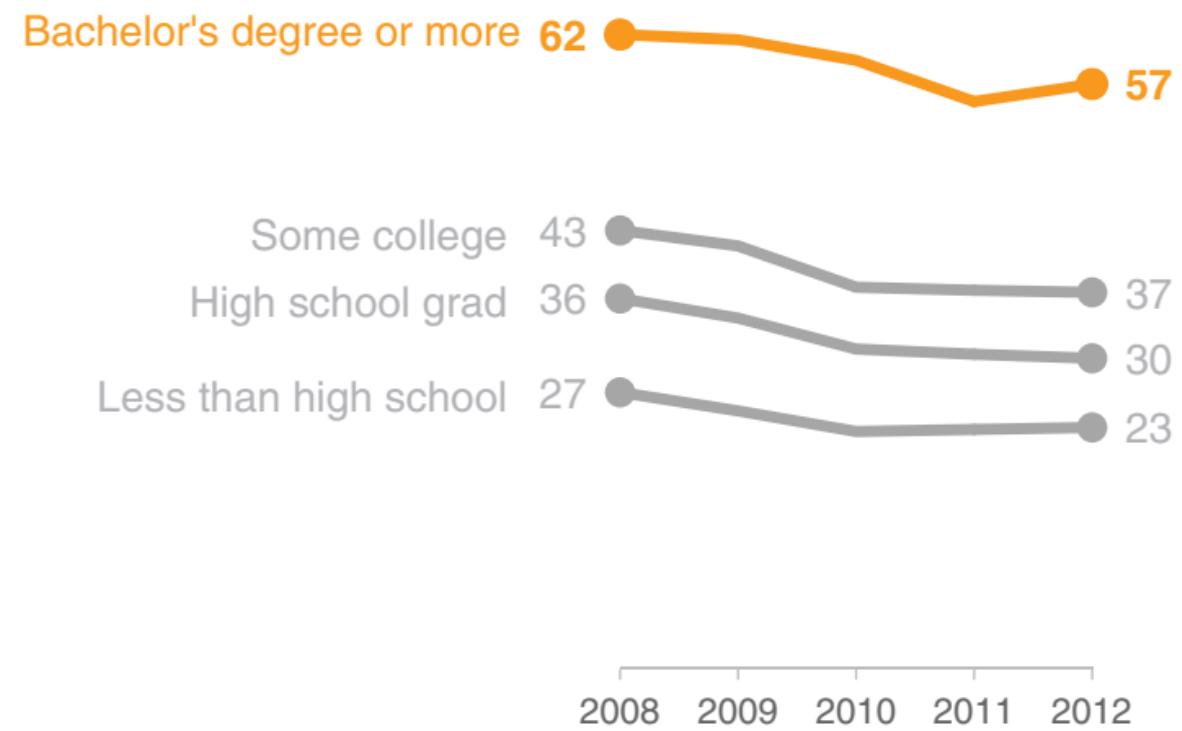
New Marriage Rate by Education

Number of newly married adults per 1,000 marriage eligible adults



New marriage rate by education

Number of newly married adults per 1,000 marriage eligible adults



Follow us on



govbigdata



Twitter



Blockdit



YouTube



Government Big Data Institute (GBDi)

Line Official



@gbdi

