

# Visualisation and data analysis for journalism studies

Rafal Urbaniak and Nikodem Lewandowski  
(University of Gdańsk)  
<https://rfl-urbaniak.github.io/teaching/>  
[rfl.urbaniak+teaching@gmail.com](mailto:rfl.urbaniak+teaching@gmail.com)

# The plan

Motivations, goals, game rules

Some history

The role of perception

Getting started with R, RStudio and ggplot2

More on what to show

Focus

Epistemic problems

Technical and mathematical problems

Statistical learning and probabilistic thinking

Statistical and analytical blunders

Basics of Bayesian thinking

Linear models

Causality and variable selection

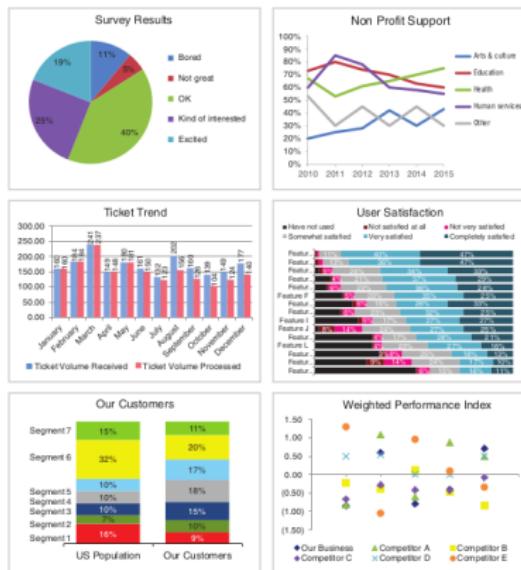
## Motivations

- It's too easy to generate tables and visualisation.
- This makes communication harder!

# Motivations

- It's too easy to generate tables and visualisation.
- This makes communication harder!

Bad graphs everywhere!



## Lack of background

- We learn some math at school.
- We learn some arts at school.

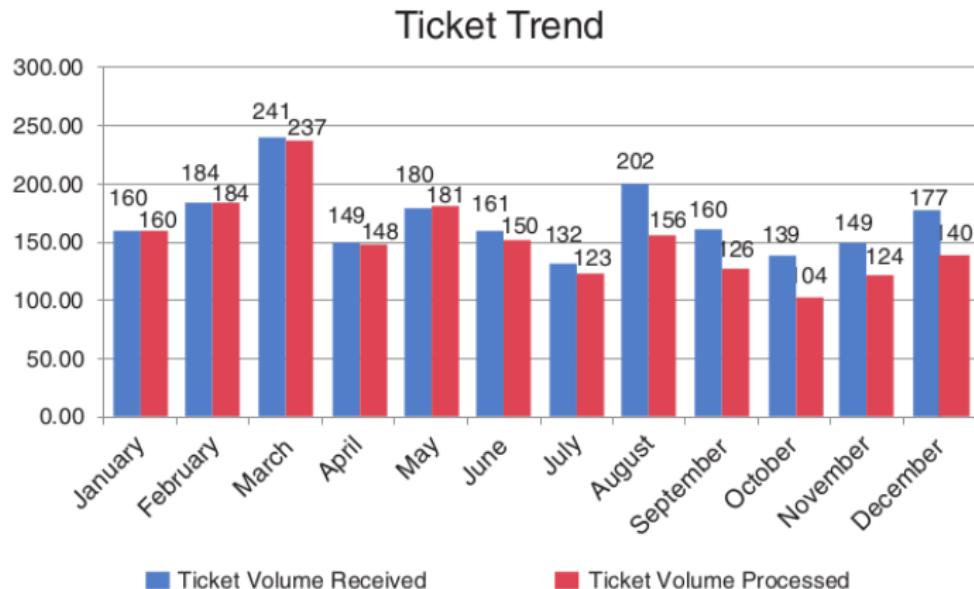
## Lack of background

- We learn some math at school.
- We learn some arts at school.

### Problem

We never learn to put them together, and think they're opposite.

## Some examples



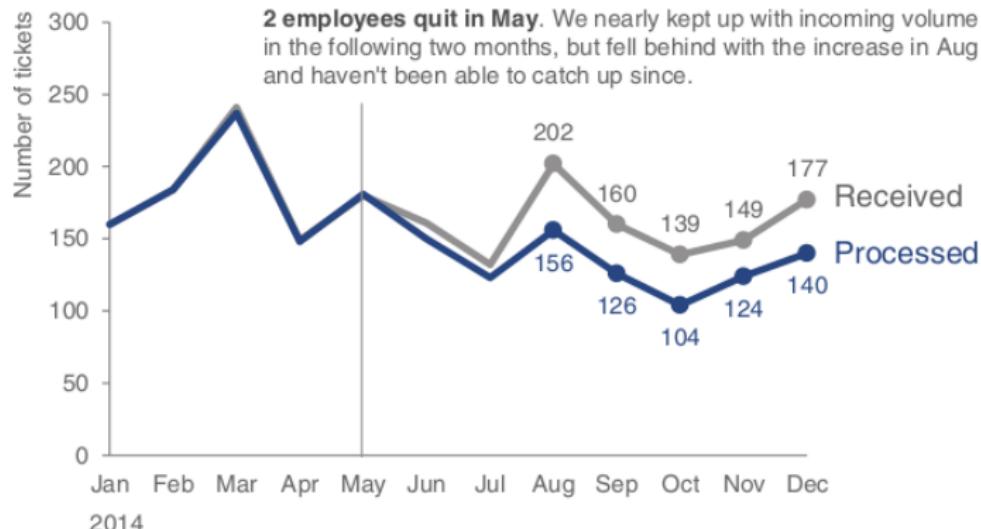
Cole Nussbaum [4]

## Some examples

### Please approve the hire of 2 FTEs

to backfill those who quit in the past year

Ticket volume over time



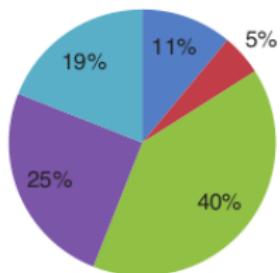
Data source: XYZ Dashboard, as of 12/31/2014 | A detailed analysis on tickets processed per person and time to resolve issues was undertaken to inform this request and can be provided if needed.

# Some examples

## Survey Results

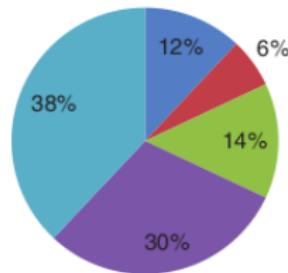
PRE: How do you feel about doing science?

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



POST: How do you feel about doing science?

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



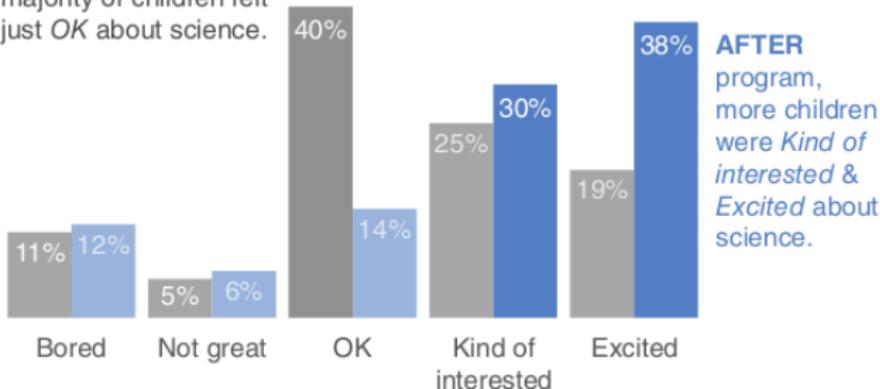
Cole Nussbaum [5]

# Some examples

## Pilot program was a success

How do you feel about science?

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

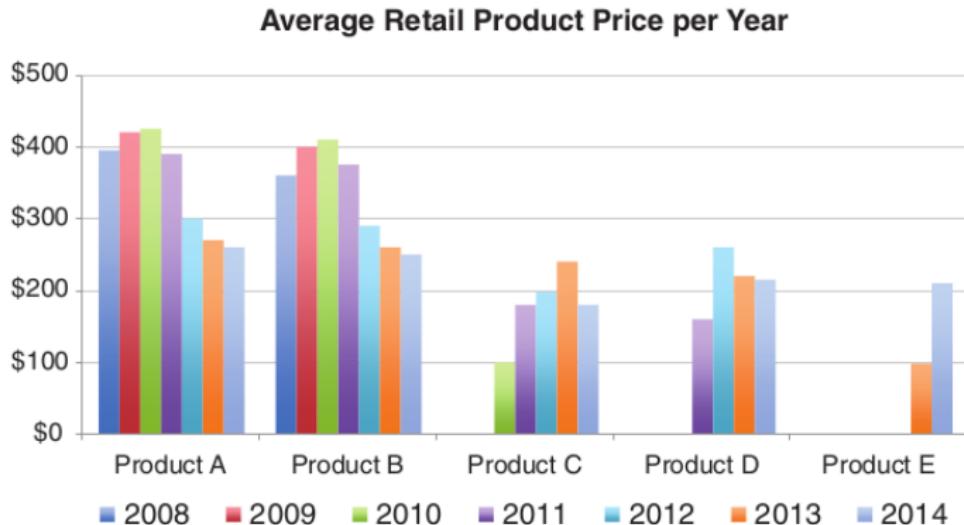


**AFTER** program,  
more children  
were *Kind of  
interested &  
Excited* about  
science.

Based on survey of 100 students conducted before and after pilot program (100% response rate on both surveys).

Cole Nussbaum [5]

## Some examples

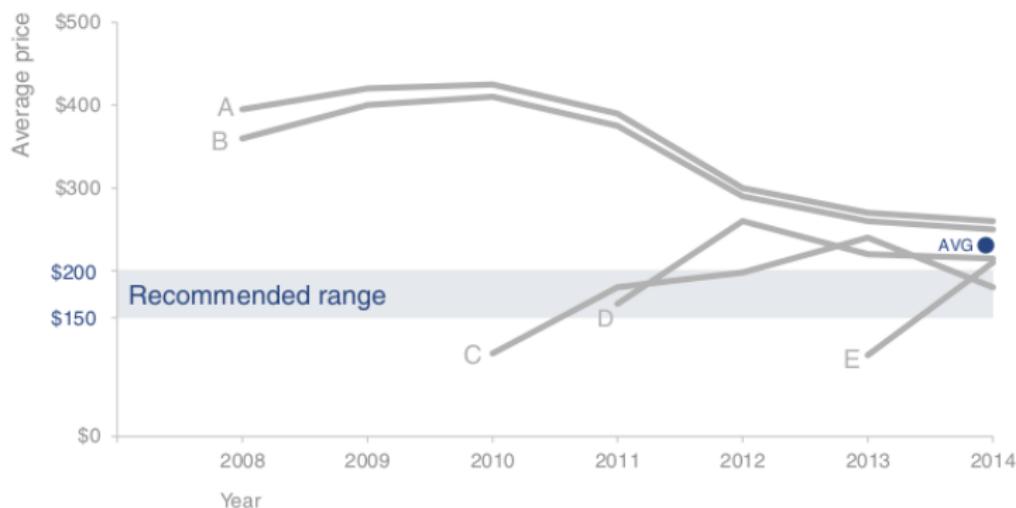


Cole Nussbaum [6]

## Some examples

To be competitive, we recommend introducing our product *below* the \$223 average price point in the **\$150–\$200 range**

Retail price over time by product



Cole Nussbaum [6]

# Goal

- To understand psychological factors that guide various visualization choices
- To be able to properly analyze data yourself (at a decent level, or at least to understand some of the complexities involved)
- To be able to visualize your data insights so that they clearly convey your message
- To be able to work in R, a statistical programming language

## Rules: final grade

### Final test: 60 points (optional)

- multiple choice with penalty points

### Project: 60 points (optional)

- two-three pages of meaningful text with at least two visualizations, bonus points for animations
- everything prepared in R markdown
- feedback loop: idea -> draft -> feedback -> revisions -> f2 -> r2

### Tutorial performance: 60 points (optional)

- If you complete a free-fall exercise without much help, show us, get some points!

## Final grade

As if out of 100.

# Contact

Updates - only here!

<https://rfl-urbaniak.github.io/teaching/>

Contact - only here!

rfl.urbaniak+teaching@gmail.com

# Sources



## Avoiding Data Pitfalls

How to Steer Clear of Common Blunders  
When Working with Data and Presenting  
Analysis and Visualizations

Ben Jones

WILEY

# Sources



cole nussbaumer knaflic

## storytelling with **data**

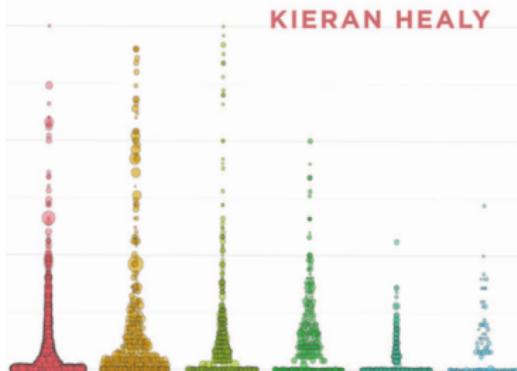
a data  
visualization  
guide for  
business  
professionals

WILEY

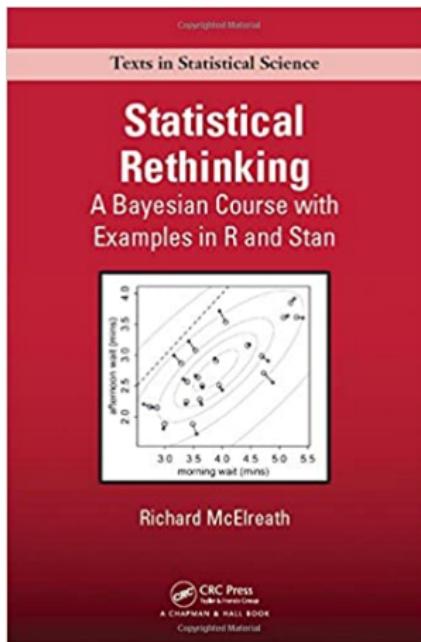
## DATA VISUALIZATION

A PRACTICAL INTRODUCTION

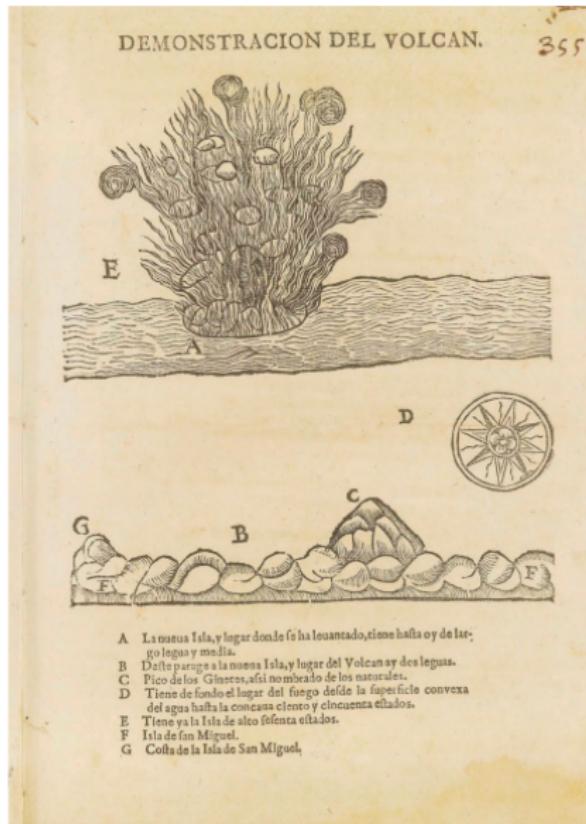
KIERAN HEALY



# Sources



# Precursors



# Precursors

B

Numb. 116.

## The Daily Courant.

Saturday, September 12. 1702.

LONDON, Sep. 12.

**T**HIS before the Duke of Ormond is  
in the Bay of Cadiz, and makes  
vaine Excuse that the English have  
undertaken a War for this last Blooded  
King. His Majesties Forces have  
got in their place a sufficient Number of Men  
and Ships to make a compleat Invasion of the Island  
and City, which may be of force off for the strict  
guarding of the Harbour that is come already  
into the hands of the French. And soe we have  
nothing else to doe but to seeke out of the last Foreign Prints  
and News-Letters, even wh' we had no already  
given what a soft measured it is.

**T**he Island of Cadiz lies between the Mouth  
of the River Guadalquivir, and the Mouths of  
the Rivers of Alfonso and Segura. It is  
joined to the Continent by a narrow Neck  
to the first Land by a Bridge of two Paes long,  
call'd le Paes de Sasse. The Distance between Paes

St. Mary and that Bridge is about 12 Miles; from  
that Bridge to the City of Cadiz is likewise about  
12 Miles: The Rocks call'd the Diamond and Los  
Picos make a Returne into the Bay pretty Dang-  
erous. The Bay is a very large one, and  
is defend'd by a small part (one Side  
which is defend'd by a Part call'd the Point, and  
the other by a Part call'd the Head) of a Mile  
and an half over: The Part of the Head on  
which the Town stands is defend'd towards the Sea  
by the Point of the Town, and towards the  
Walls of the Town (which meet the Part are  
walk'd by the Tide), and by Sharp Rocks: And there  
are very strong Fortifications to secure the Passage  
over the Town, and the Mouth of Land, which  
rises from the larger Part of the Island in the City.  
The Town of St. Catherine which is taken by the  
French, is a very strong Place, and well built  
There are a great many Churches in Cadiz, which  
is well built, very rich, and full of Inhabitants.

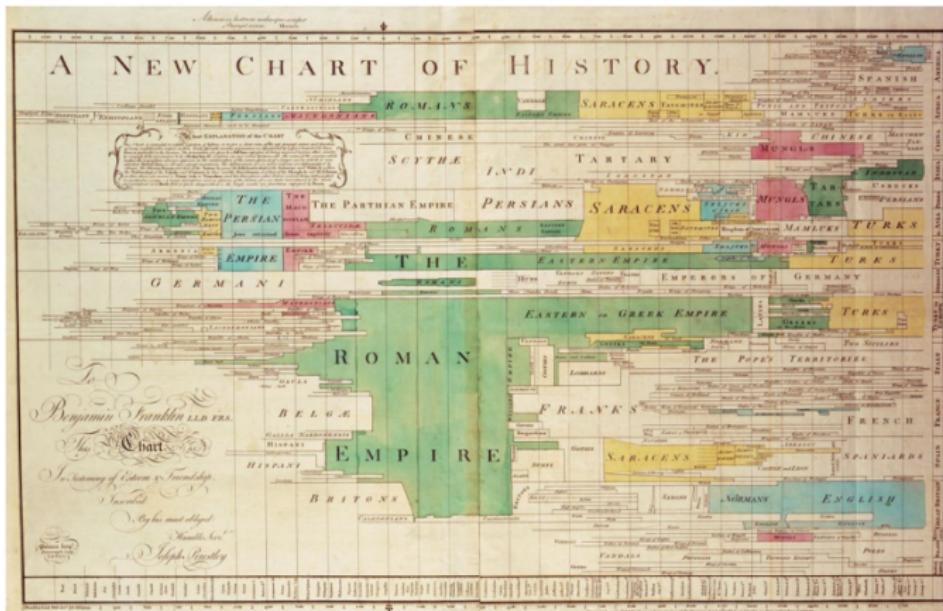


A. A. The Bay of Cadiz.  
B. The City.  
C. The Harbour.  
D. The Rocks.  
E. Pier St. Mary.  
F. Pier Rock.  
G. Pier Sasse de Sasse.  
H. Pier of the Señor.  
I. The fortifications on the Neck of Land be-  
tween the City.  
L. The City.  
K. The Point.  
L. Los Picos.  
M. Cadiz.  
N. Point of St. Sebastian.  
O. Point of Cadiz.  
P. The Island of St. Peter.  
Q. A Bank of two Miles.

English *The Daily Courant* (invasion on Cadiz), 1702

# Precursors

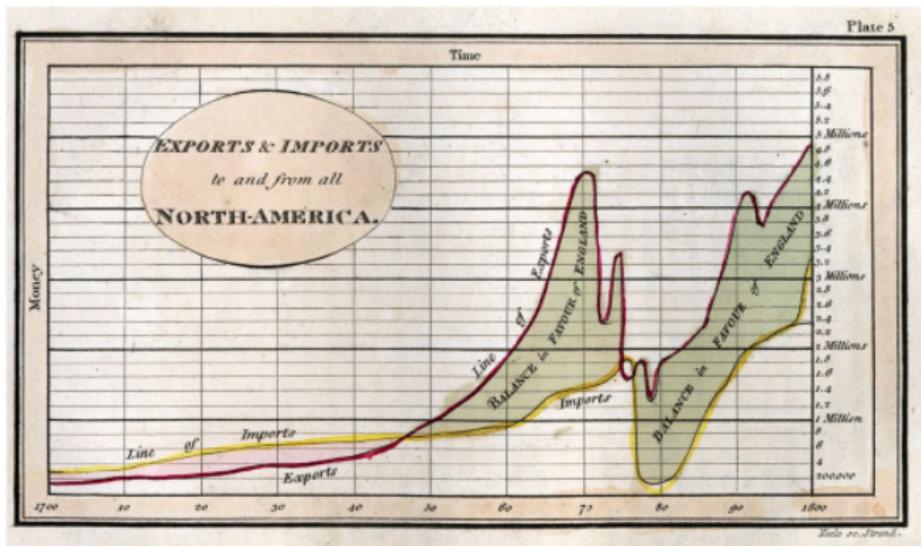
## Joseph Priestley (1733-1804)



A new chart of history, 1769

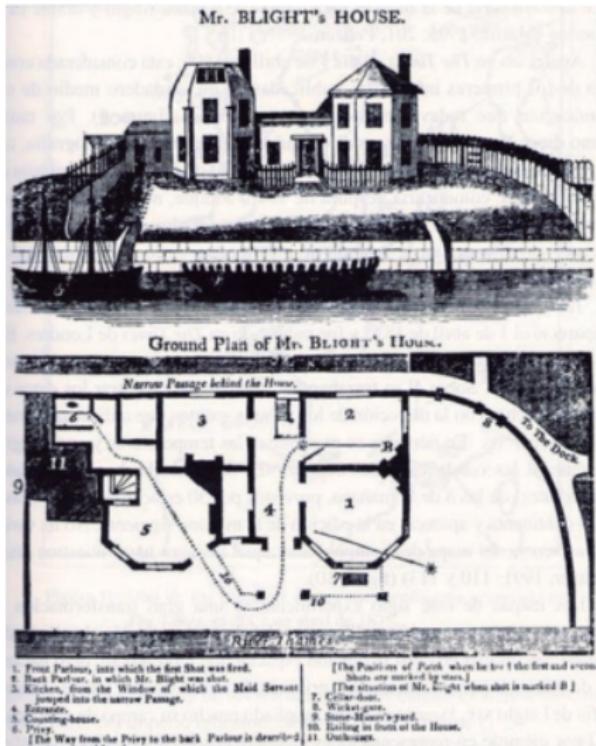
## Precursors

## William Playfair (1759-1823)



*Statistical breviary, 1801*

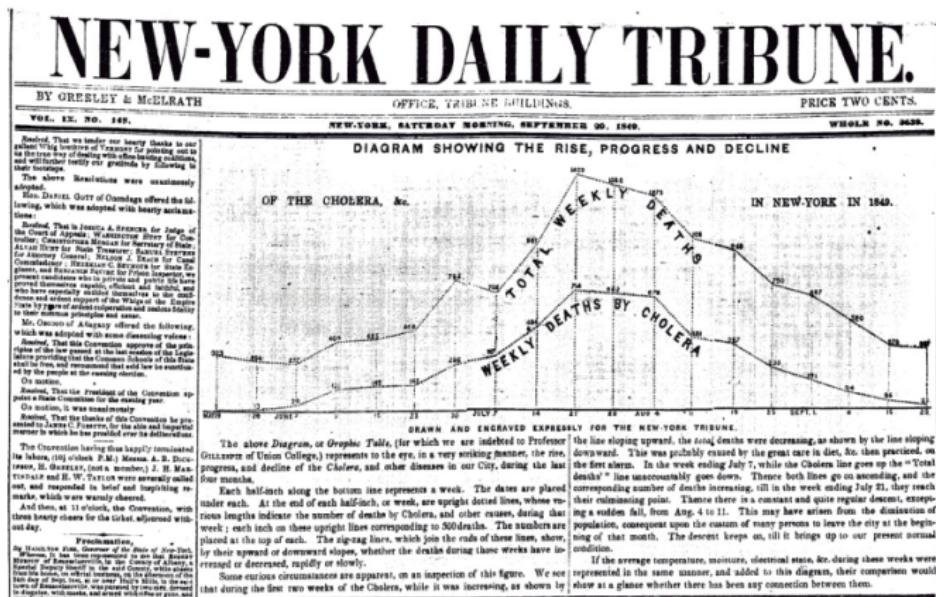
# Precursors



A murder case coverage in *The Times*, 1806

## Precursors

## William Mitchell Gillespie (1816-1868)



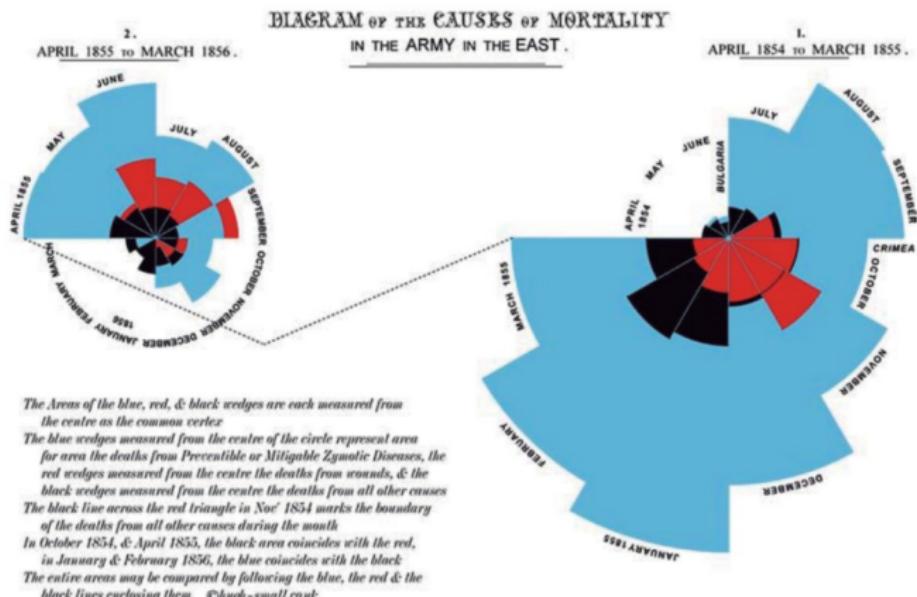
New-York Daily Tribune, 1849

# XIXth century explosion

## Reasons

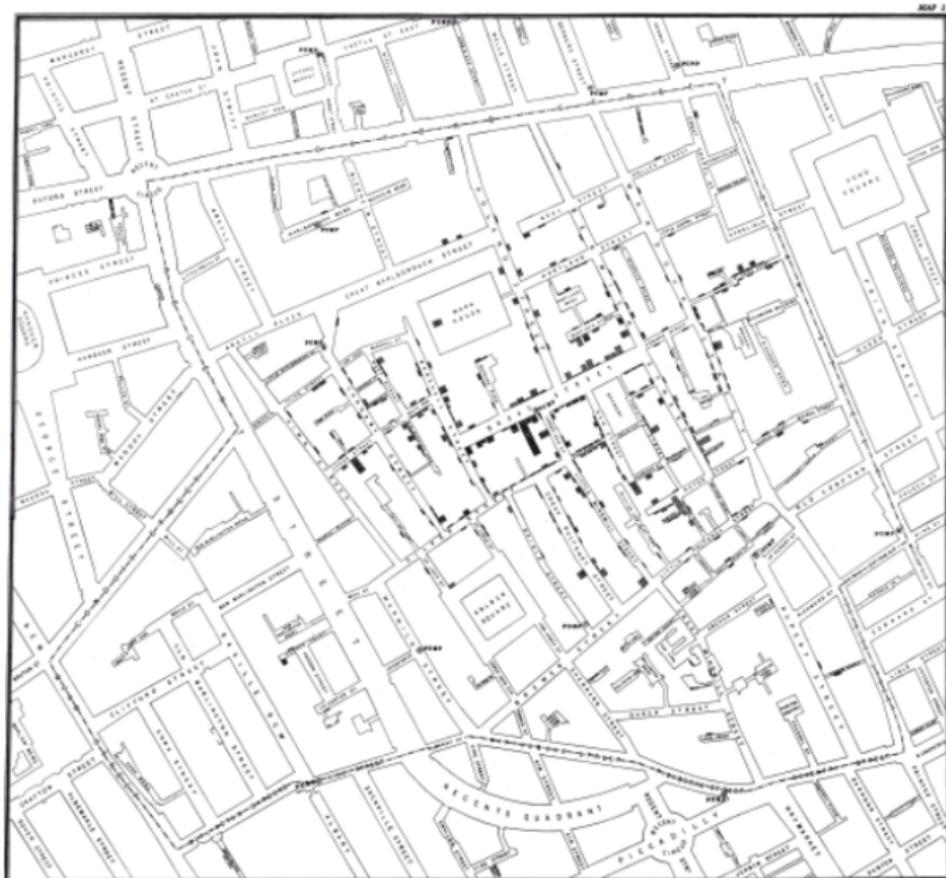
- modern nation-states with increased interest in collecting economic and demographic data
- descriptive statistical methods used before in physical sciences began to be used in social sciences (e.g. Adolphe Quetelet, Francis Galton)
- dawn of new sciences, such as epidemiology

# Florence Nightingale (1820-1910) and the Crimean war

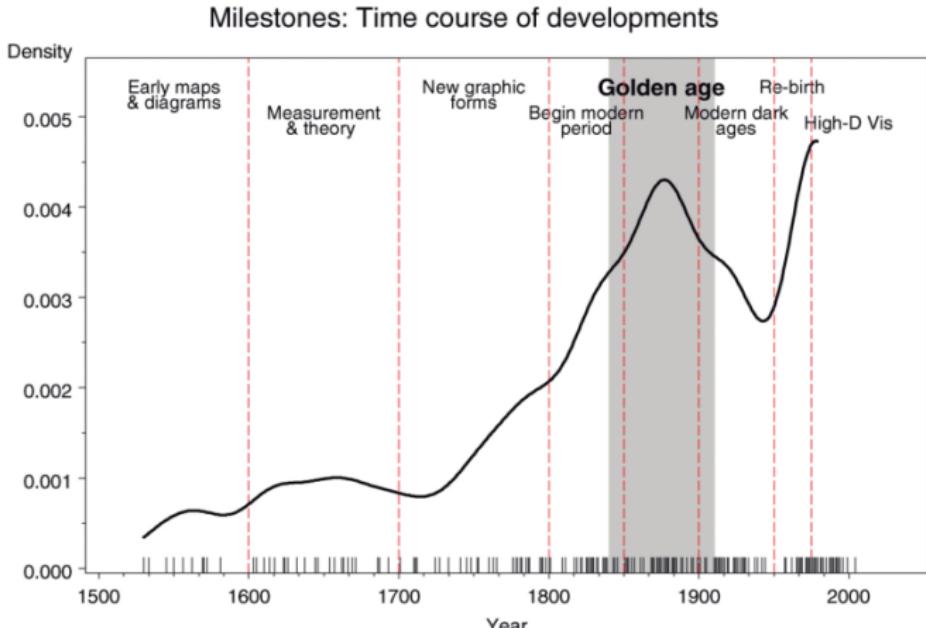


Causes of Mortality, 1856

# John Snow (1813-1858) and cholera in London



# Modern dark ages in statistics



Number of visualization historical landmarks per year, *Friendly 2008*

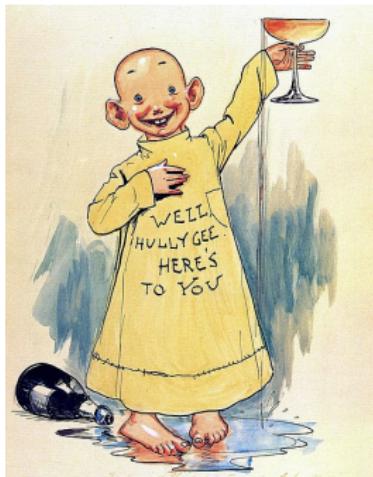
## The pictorial turn in newspapers

*Newspapers became a prime site where visual art and popular forces met and made their peace, and news contributed to the fullness of modernism as it arrived in the twentieth century [...] During the century, the newspapers in the study shifted from the abundant complexity of the Victorian era to the fixed simplicity of modernism. They adopted all the specific forms commentators identified with the modern style: fewer columns, prominent illustrations, horizontal layout, and simplified headline typography.* (Barnhurst & Nerone 2001)

# Yellow kid journalism (1895-1898)

## Say what?

Sensational journalism in the circulation war between Joseph Pulitzer's *New York World* and William Randolph Hearst's *New York Journal* (Pulitzer tried to be more content-based but circulation shrank)



Yellow Kid, *New York World* and *New York Journal*

# Yellow kid journalism (1895-1898)

**\$50,000 REWARD.—WHO DESTROYED THE MAINE?—\$50,000 REWARD.**  
EDITION FOR GREATER NEW YORK  
**NEW YORK JOURNAL**  
AND ADVERTISER

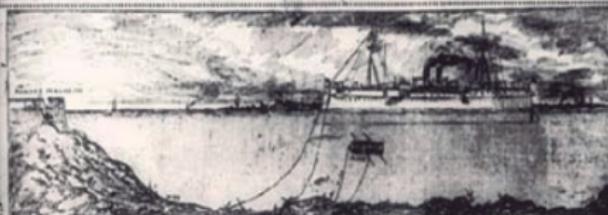
**DESTRUCTION OF THE WAR SHIP MAINE WAS THE WORK OF AN ENEMY**

**\$50,000!**  
\$50,000 REWARD!  
for the Detection of the  
Perpetrator of  
the Maine Outrage!

Assistant Secretary Roosevelt  
Convinced the Explosion of  
the War Ship Was Not  
an Accident.

The Journal Offers \$50,000 Reward for the  
Conviction of the Criminals Who Sent  
anti-American Sailors to Their Death.  
Naval Officers Unanimous That  
the Ship Was Destroyed  
on Purpose.

**\$50,000!**  
\$50,000 REWARD!  
for the Detection of the  
Perpetrator of  
the Maine Outrage!



**NAVAL OFFICERS THINK THE MAINE WAS DESTROYED BY A SPANISH MINE.**

Hidden Mine or a Sunken Torpedo Believed to Have Been the Weapons Used Against the American Men-of-War—Officers and Men Tell Thrilling Stories of Being Blown Into the Air by a Mass of Shattered Steel and Exploding Shells—Survivors Brought to Key West Now Tell Identical Accidents—Spanish Official Proves the Machado's Secret Order a Seizing Inquiry—Journal Sends Divers to Havana to Report Upon the Condition of the Wreck.

By C. OTIS SMITH, U. S. PLAINFIELD, N. J.

Illustrations of scenes of carnage and destruction, such as were brought about by the explosion of the Maine in Havana Harbor, were not available. Illustrations of the "New England Round" can be supplied, showing that the crew of the Maine in Havana Harbor was not an ordinary naval crew, but one of the best and bravest in the world.

The suggestion that the Maine was destroyed by a mine is not supported by the facts.

The suggestion that the Maine was deliberately blown up, given stronger every hour. Not a single fact to the contrary has been produced.

Captain Ingles, of the Maine, and his comrade, Captain Lee, both said that public opinion be suspended until they have completed their investigation.

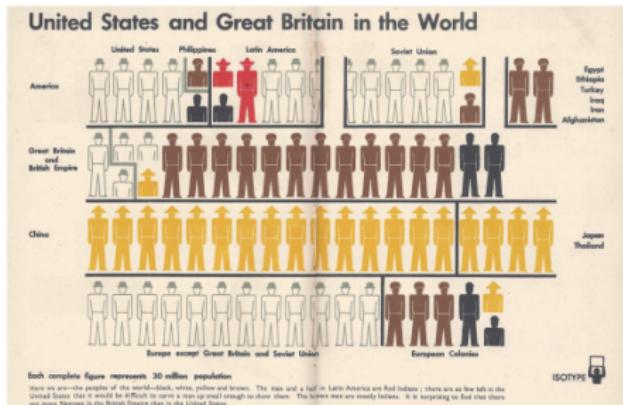
They are taking the names of naval men who are connected with their families, and publishing them.

Washington agents can tell that Captain Ingles and Captain Lee are each other as a hidden man. The English eight o'clock was used all day yesterday.

The sinking of Maine in the bay of Havana (notice the Spanish mine), *New York Journal*, Feb. 17, 1898

# Viennese Museum for Society and the Economy (1924)

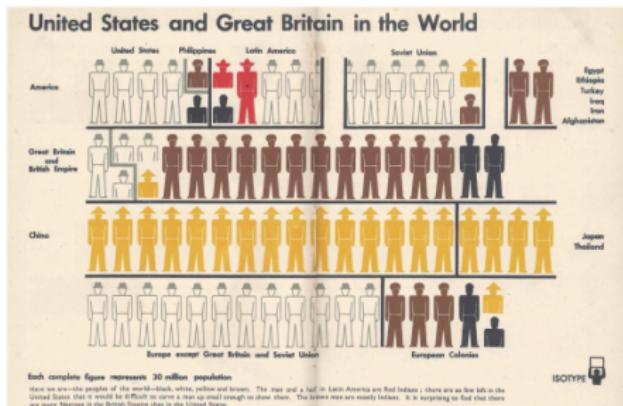
## Facts for the uneducated



ISOTYPE, universal visual language by Neurath, Arntz and Reidemeister

# Viennese Museum for Society and the Economy (1924)

## Facts for the uneducated



ISOTYPE, universal visual language by Neurath, Arntz and Reidemeister

## The “Bible”

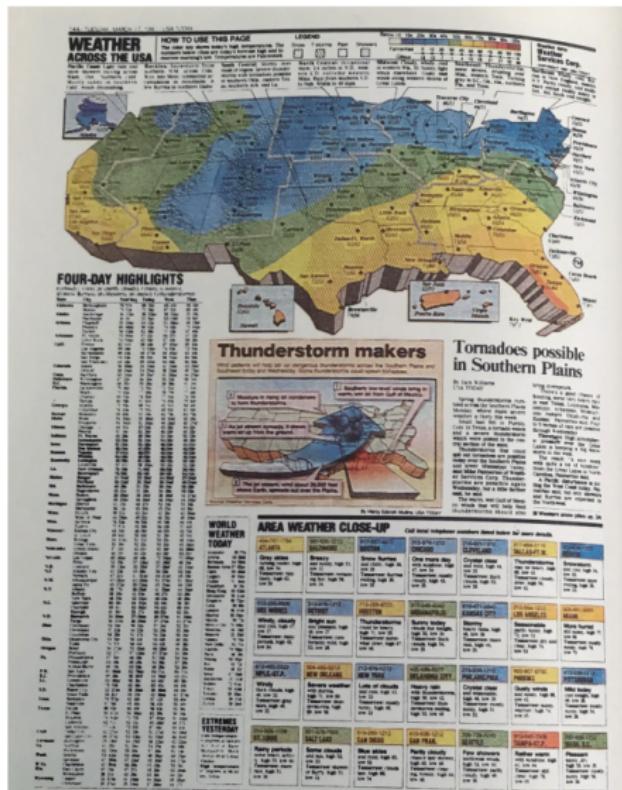
Pictographs and Graphs: How to Make and Use Them, Modley & Lowenstein, 1952

# ISOTYPE



A page from *Fortune*, 1929

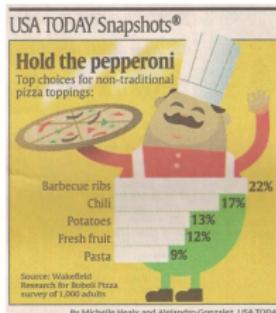
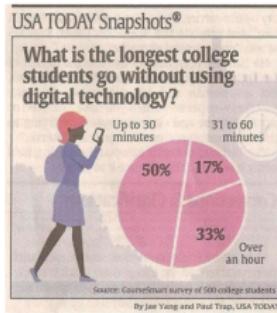
# Birth of USA Today (1982)



A revolutionary weather map

# Birth of USA Today (1982)

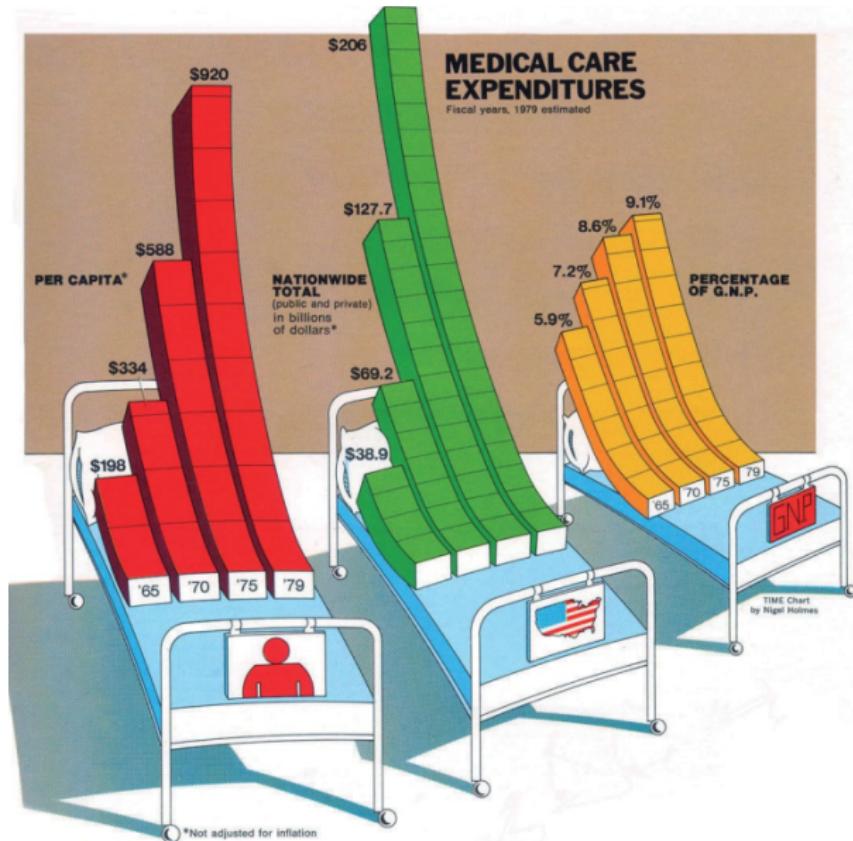
- its success expanded the use of graphics in print publications
- tilted the stylistic balance towards the pictorial and lighthearted
- art training, no quantitative expertise
- in 1984 60% of 156 newspapers reported an increased use of news graphics, and an additional 22% said that they had just incorporated them into their pages



## What's the problem?

*Nearly all those who produce graphics for mass publication are trained exclusively in the fine arts and have had little experience with the analysis of data [...] Illustrators too often see their work as a exclusively artistic enterprise—the words "creative", "concept", and "style" combine regularly in all possible permutations, a Big Think jargon for the small task of constructing a time-series a few data points long. Those who get ahead are those who beautify data, never mind statistical integrity.*

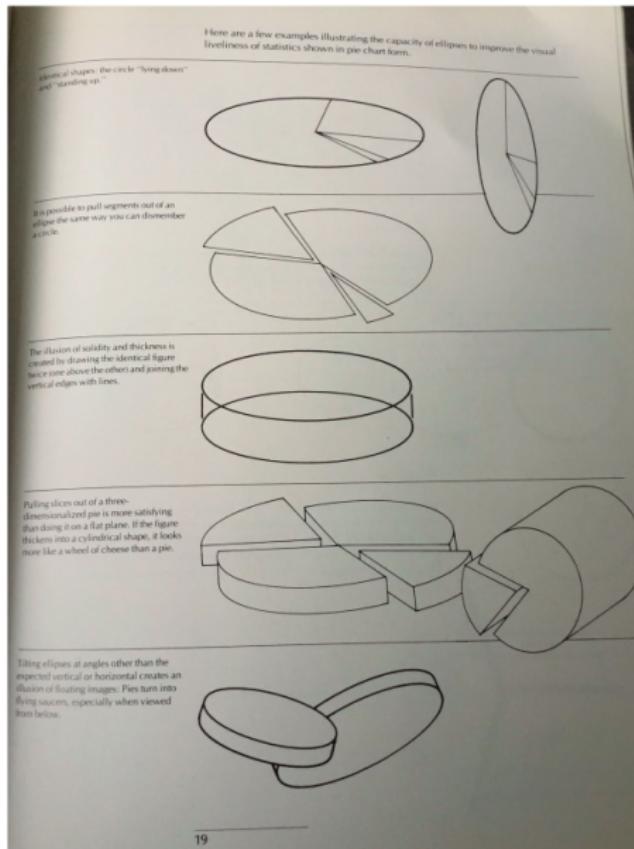
*[Edward Tufte 1983]*



# Nigel Holmes

*As long as the artist understands that the primary function is to convey statistics and respect that duty, then you can have fun (or be serious) with the image: that is, the form in which those statistics appear. Boredom is as much a threat in visual design as it is elsewhere in art and communication. The mind and eye demand stimulation and surprise.*

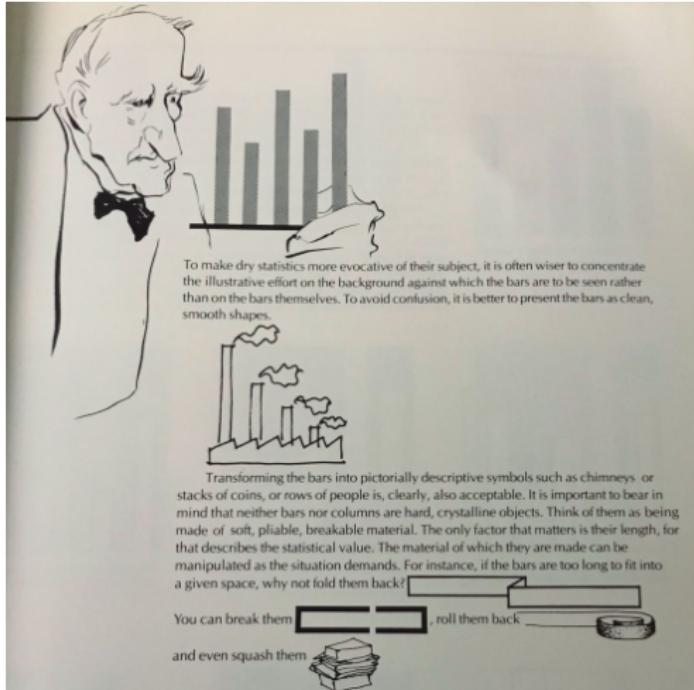
# Jan V. White



*To make dry statistics more evocative of their subject, it is often wiser to concentrate the illustrative effort on the background against which the bars are to be seen rather than on the bars themselves, [...] transforming the bars into pictorially descriptive symbols such as chimneys or stacks or coins, or rows of people is, clearly, also acceptable [...] The material of which they are made can be manipulated as the situation demands. For instance, if the bars are too long to fit into a given space, why not fold them back? You can break them, roll them back and even squash them.*

*(Jan. V. White, 1984)*

Jan V. White



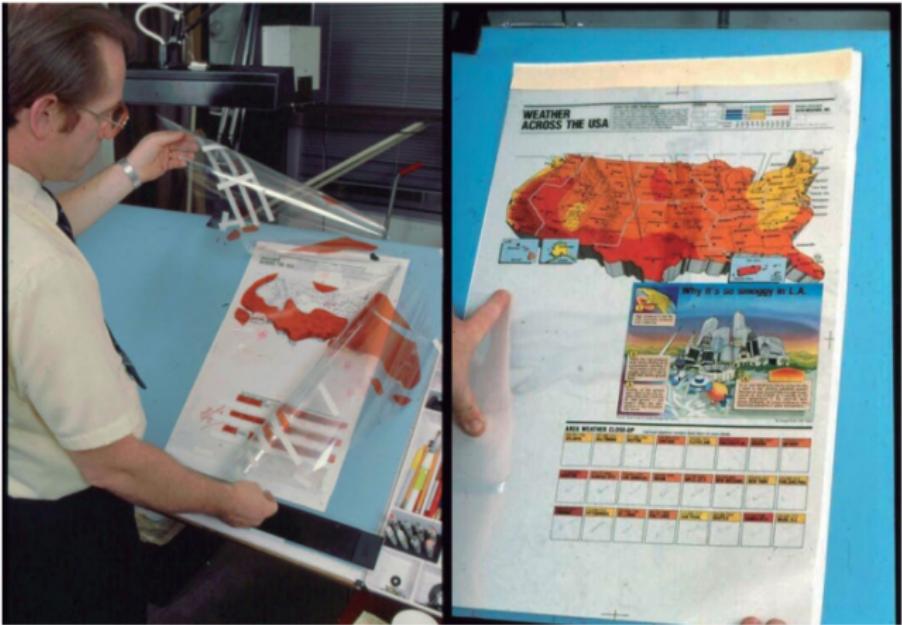
White's textbook on visualization, 1984

# Computer-age graphics



George Rorick, hand-made visualisation, 11 a.m. to 6 p.m.

# Computer-age graphics



George Rorick, hand-made visualisation, 11 a.m. to 6 p.m.

## Computer-age graphics

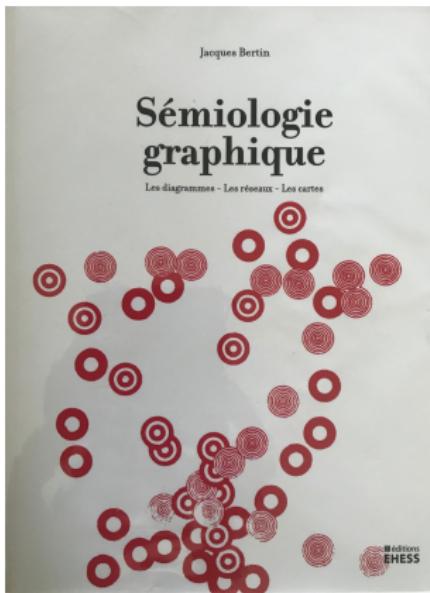
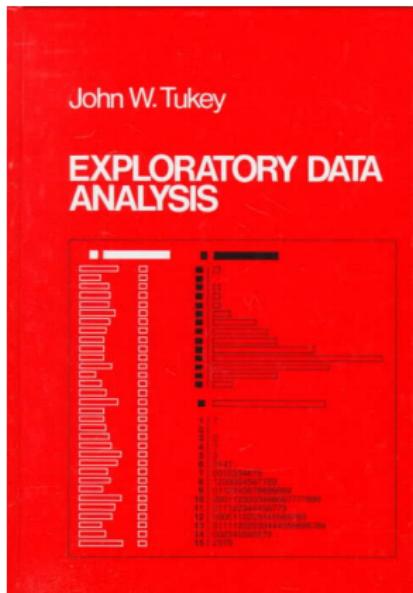
- Apple, 1984
- PostScript & Adobe Illustrator, 1987 (raster vs. vector files)
- Adobe Photoshop, 1989

*We went from some very nice illustrated graphics to some very poor computer-generated graphics, but that was the limitations of the technology, and it took about at least five years, maybe more, before we started to see the computer graphics start to rise up in quality.*

*John Grimwade (check out his website!)*

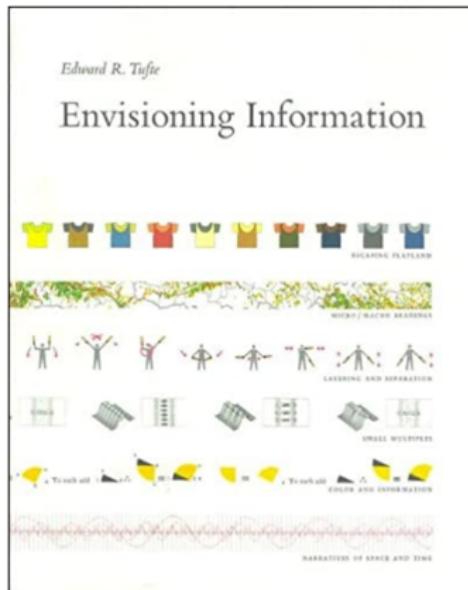
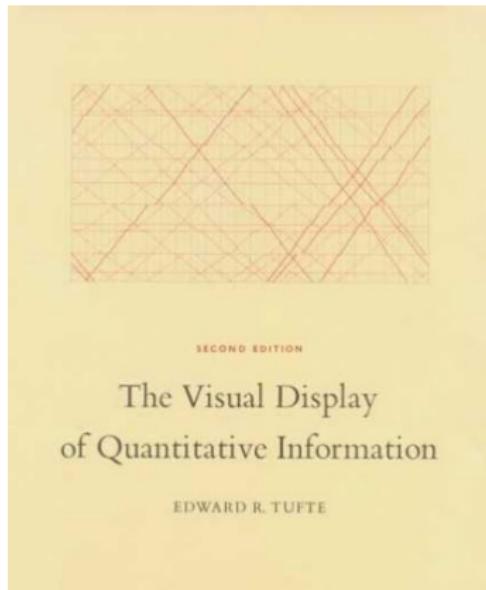
# Backlash against cartoons

Tukey 1977, Bertin 1967



# Backlash against chartoons

Tufte 1983, 1990



## Backlash against cartoons

*Sometimes decoration can help editorialize about the substance of the graphic. But it is wrong to distort the data measures —the ink locating values of numbers—in order to make an editorial comment or fit a decorative scheme.*

*(Tufte 1983: 59)*

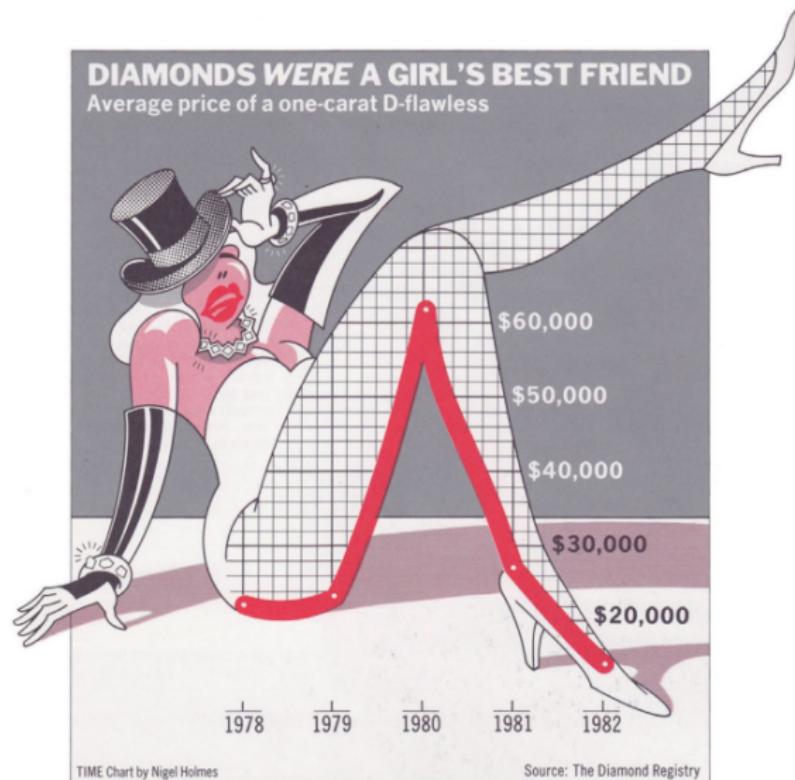
## Backlash against cartoons

*If you belong to the school of people who believe that charts should only present statistics in the most straightforward, plain way, with no other visual help to the reader, for example, than the bar of the bar chart, the line of the fever graph, the circle of the pie chart, or the rules of the table, then move on to another part of the book [...] Boredom is as much a threat in visual design as it is elsewhere in art and communication. The mind and eye demand stimulation and surprise [...] Even a smile will encourage a reader to look into the statistics he or she might not have thought of reading in a less embellished chart.* (Holmes 1984: 72)

## Backlash against cartoons

*Too many data presentations [...] seek to attract and divert attention by means of display apparatus and ornament. Chartjunk has come to corrupt all sorts of information exhibits and computer interfaces (Tufte 1990: 33)*

# Backlash against cartoons



Holmes' chart in the *Times* magazine

## Backlash against chartoons

*Consider this unsavory exhibit at right —chockablock with cliché and stereotype, coarse humor, and a content-empty third dimension. Is it the product of a visual sensitivity in which a thigh-graph with a fishnet-stocking grid counts as Creative Concept. [...] Lurking behind chartjunk is contempt for both information and for the audience. Chartjunk promoters imagine that numbers and details are boring, dull, and tedious, requiring ornament to enliven. Cosmetic decoration, which frequently distorts the data, will never salvage an underlying lack of content. If the numbers are boring, then you've got the wrong numbers. Credibility vanishes in clouds of chartjunk; who would trust a chart that looks like a video game? (Tufte 1990: 34).*

## Backlash against cartoons

*Graphical competence demands three quite different skills: the substantive, statistical, and artistic. Yet now [in the early 80s] most graphical work, particularly at news publications, is under the direction of but a single expertise —the artistic. Allowing artist-illustrators to control the design and content of statistical graphics is almost like allowing typographers to control the content, style, and editing of prose.*  
*(Tufte 1983: 87).*

## Recent developments

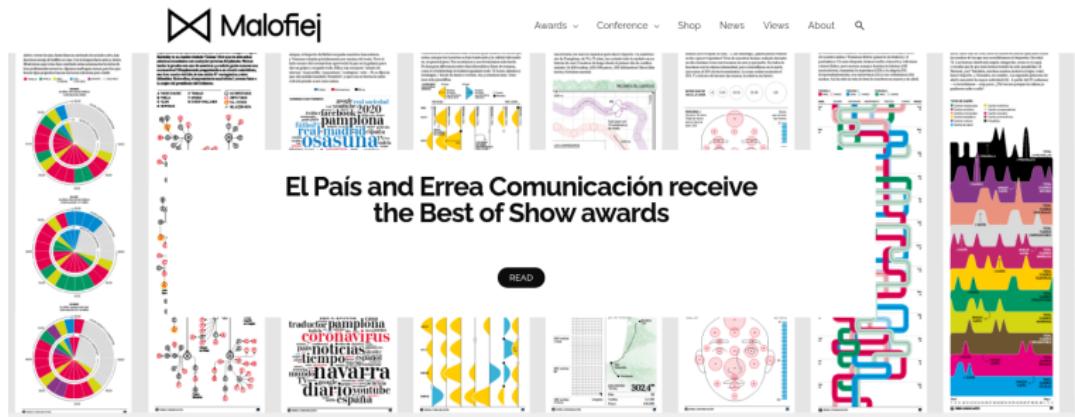
# Recent developments

## Geek takeover

- more information density and more data
- visualization desks more independent from arts departments
- the 90s and early 2000s: illustration-driven explanations, sometimes supplemented by small and straight-forward statistical graphs and data maps
- today, the balance has shifted to presentations that rely mainly on the visual display of data, both quantitative and qualitative
- often, no longer detached “graphics departments”. Data journalists, nerd journalism!

# Recent developments

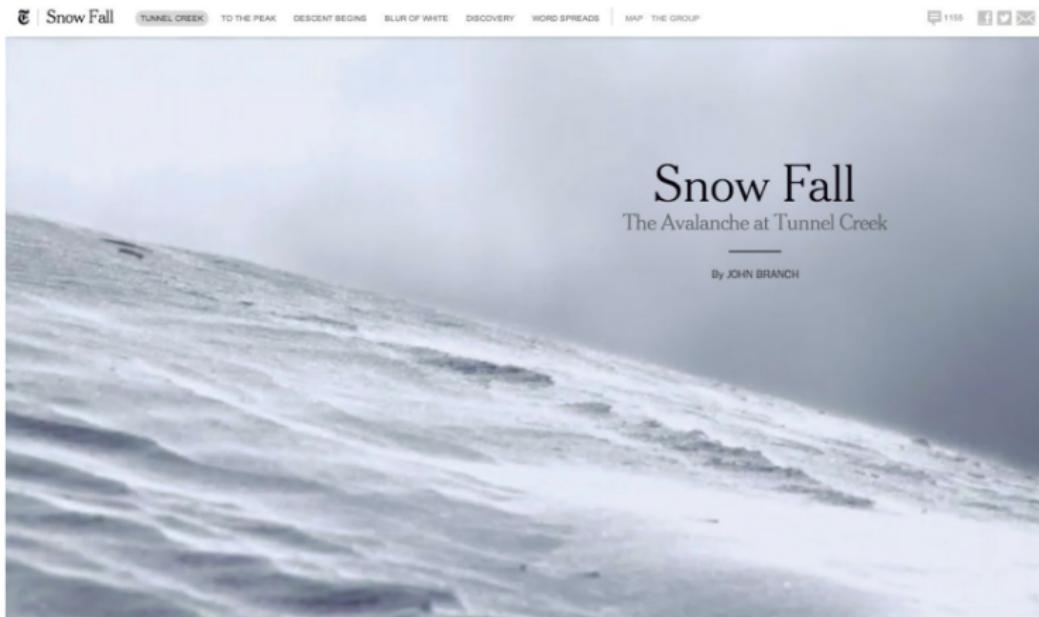
Check out Malofiej awards (1992)



Malofiej awards website

# Recent developments

Example (“new era”, 3 mln. in no time)



Snowfall at NY Times

# Recent developments

## Example (most popular piece in Times, 2013)

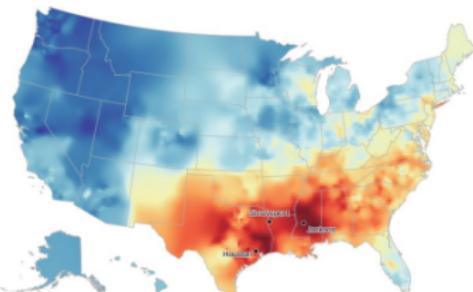
### How Y'all, Youse and You Guys Talk

What does the way you speak say about where you're from?  
Answer all the questions below to see your personal dialect map.

#### Your Map

See the pattern of your dialect in the map below. Three of the most similar cities are shown.

Least similar   Most similar   Show least similar   SHARE YOUR MAP:



These maps show your most distinctive answer for each of these cities.

JACKSON



What do you call a drive-through liquor store?

beer barn

HOUSTON



What do you call the small road parallel to the highway?

feeder road

SHREVEPORT



What do you call a sweetened carbonated beverage?

coke

How Y'all quiz, NYT

## For the tutorial

Complete the introductory instructions about github, bring a flash drive!