

INST0065

Data Visualization and GIS

Week 2: Data types and visualization types, and, a history of visualization

Dr. Oliver Duke-Williams

o.duke-williams@ucl.ac.uk

(Please use Moodle forums for messages about this module)

Twitter: @oliver_dw



INFORMATION STUDIES

Celebrating 100 years

Contents

- Types of data
- Types of visualization
- A brief history of data visualization

Types of data, types of story

- This section based on:

Few, S. (2004). *Eenie, Meenie, Minie, Moe: Selecting the Right Graph for Your Message.*

http://www.perceptualedge.com/articles/ie/the_right_graph.pdf

Data types

- In order to plan a visualization, we need to recognise different types of data
- The primary division is between quantitative and categorical data
 - Quantitative data are numerical
 - Categorical data have labels (sometimes partially numeric)

Examples

Quantitative	Categorical
Height	Nationality
Weight	Day of the week
Age	Whether vowel/consonant
Temperature	Gender
Price	Parts of speech

Sub-divisions of categorical data

- Nominal
 - Categories have labels but have no inherent order
(gender, lemma (linguistics))
- Ordinal
 - Categories have an intrinsic order
(day of the week, gold/silver/bronze medals)
- Interval
 - Categories converted from numerical data
 - (age 0-4, 5-9, 10-14...)

It is not always clear how to define these

- Colour – usually nominal, but consider wavelength spectra
- Days, months etc – usually ordinal, but need definition: when does the week start?
 - For data display: finite set
 - For data analysis: may be cyclical (Jan, Feb,...,Nov,Dec,Jan)

Dichotomous data

- Dichotomous data are nominal data that can have exactly two possible states e.g. agree/disagree
 - Often framed as: 1 – has characteristic; 0 – does not have characteristic
 - Nominal data with >2 categories might be presented for analysis as a set of dichotomous variables
 - Sometimes called binary data or binomial data; these terms are also used for other things

Intervals should be carefully considered

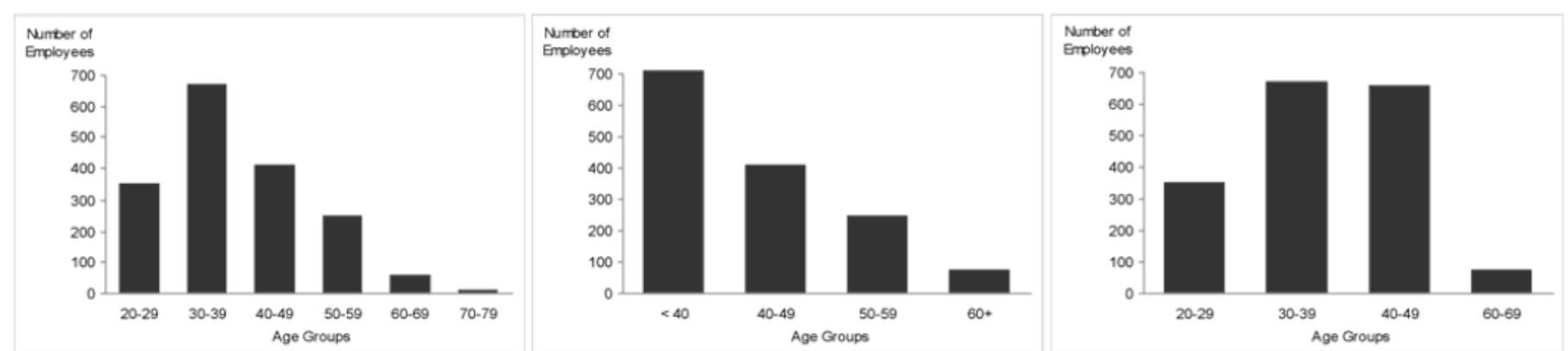


Figure 2: Changing interval scales can alter the data and result in misinformation.

20-29 30-39 40-49 50-59 60-69 70-79
Age Groups

< 40 40-49 50-59 60+
Age Groups

20-29 30-39 40-49 60-69
Age Groups

Sub-divisions of quantitative data

- Continuous vs discrete
 - Discrete data have fixed units; e.g. persons
 - Continuous data can be divided into smaller units: e.g. time

Sub-divisions of quantitative data

- Interval v ratio
 - Interval data: differences between units are the same
 - e.g. temperature
 - Ratio data is interval data with a natural zero
 - e.g. time, distance
 - Temperature is not ratio data: 20°C is not twice as hot as 10°C (consider what happens if we quote temperature in $^{\circ}\text{F}$)

Individual vs aggregate observations

- Individual records will typically have both categorical and quantitative variables
- As we group these, we will introduce the fundamental quantitative variable 'count' or 'number'

Visual attributes

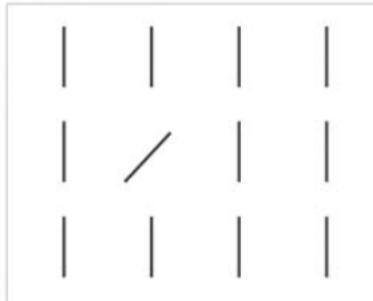
- Few (2004) identifies a set of visual attributes
 - 2-D location
 - Line length
 - Size
 - Shape
 - Orientation
 - Colour

Visual attributes

- A slightly fuller set are described in another paper (Few, 2004b)
 - Form
 - Orientation, line length, line width, size, shape, curvature, added marks, enclosure
 - Colour
 - Intensity, hue
 - Position
 - x,y

Form

Orientation



Line Length



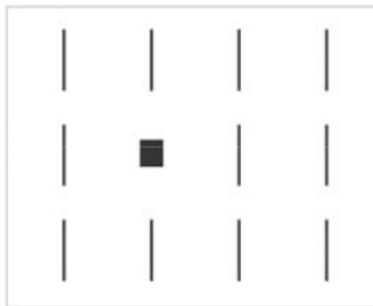
Line Width



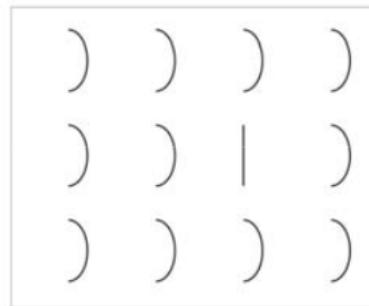
Size



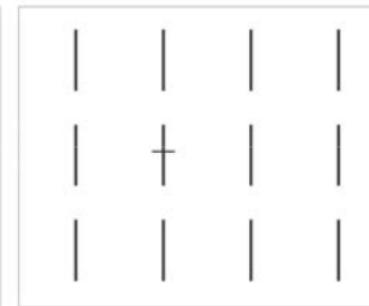
Shape



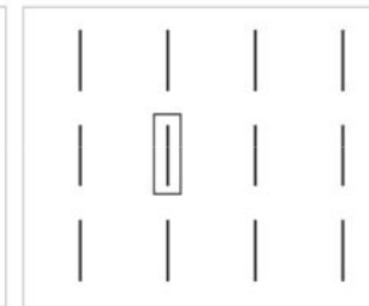
Curvature



Added Marks

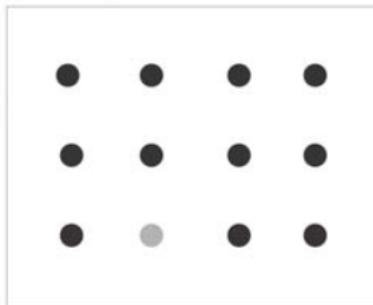


Enclosure

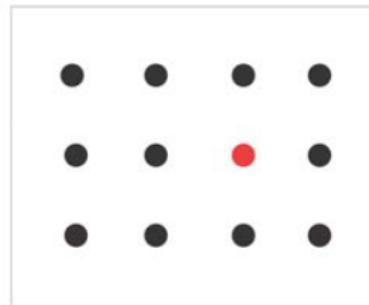


Color

Intensity



Hue



Spatial Position

2-D Position



Should we be limited to visual attributes?

- Data 'visualization' tautologically promotes visual presentation
- Data can be presented in other ways
 - Data sonification
 - Haptic data

Visual attributes

- Of these attributes, Few states:

"Of the full list of these attributes, only two emerge as highly effective means to visually encode quantitative values: 2-D location and line length."

Do you agree?

Orientation

- Orientation, as one example, might be seen as either:
 - a distinct characteristic
 - a line formed by a pair of 2D locations

- Value is shown by orientation as well as colour

COVID-19 weekly change

Small areas (MSOAs) in England



The size of the arrow shows the size of the increase or decrease (absolute change, not percentage change)

Week ending 6th January

On this map, each arrow represents a Middle-Layer Super Output Area (MSOA) with a population of around 7,000-10,000 people.

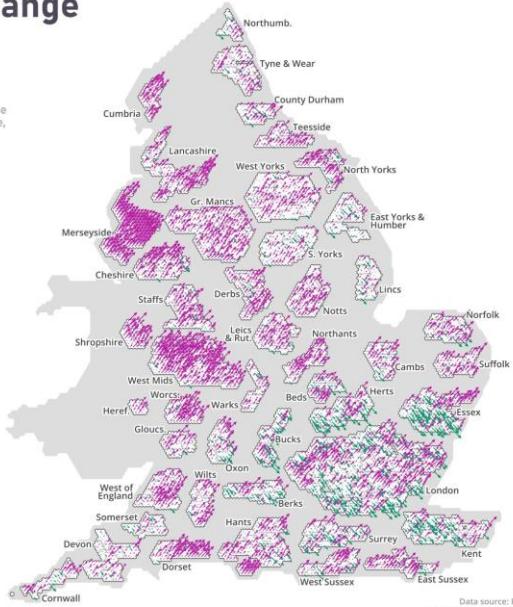
The size of the arrow shows the size of the absolute change in the COVID-19 rate (people testing positive per 100,000 population) in the week ending 6th Jan compared with the previous week.

Areas are grouped by ceremonial counties and other recognisable sub-national areas. These groups include unitary authorities (e.g. Nottingham UA in the Notts group) and don't all reflect current local gov structures.

Dark lines between adjacent areas represent local authority boundaries. Faint local authority names can be seen on the high-res version.

Grey areas between local authority groups don't represent data and serve only as a background to aid with interpreting the map.

Changes in rates can reflect changes in the numbers of tests performed as well as changes in the prevalence of the virus.



Map by @carlbaker

Data source: Public Health England

Template: tinyurl.com/HCL-hex-cartograms

Source: <https://twitter.com/carlbuster/status/1348933992100188162>

COVID-19 weekly change

Small areas (MSOAs) in England



The size of the arrow shows the size of the increase or decrease (absolute change, not percentage change)

Week ending 6th January

On this map, each arrow represents a Middle-Layer Super Output Area (**MSOA**) with a population of around 7,000-10,000 people.

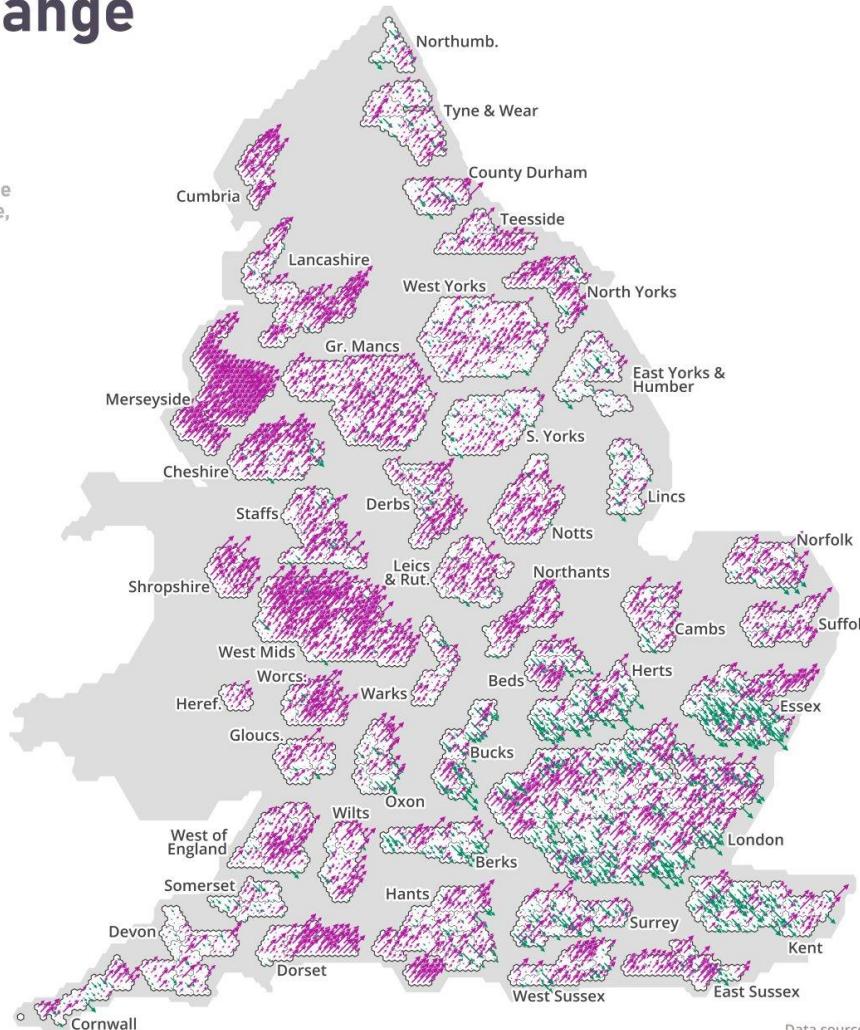
The size of the arrow shows the size of the absolute change in the COVID-19 rate (people testing positive per 100,000 population) in the week ending 6th Jan compared with the previous week.

Areas are grouped by ceremonial counties and other recognisable sub-national areas. These groups include unitary authorities (e.g. Nottingham UA in the Notts group) and don't all reflect current local gov structures.

Dark lines between adjacent areas represent local authority boundaries. Faint local authority names can be seen on the high-res version.

Grey areas between local authority groups don't represent data and serve only as a background to aid with interpreting the map.

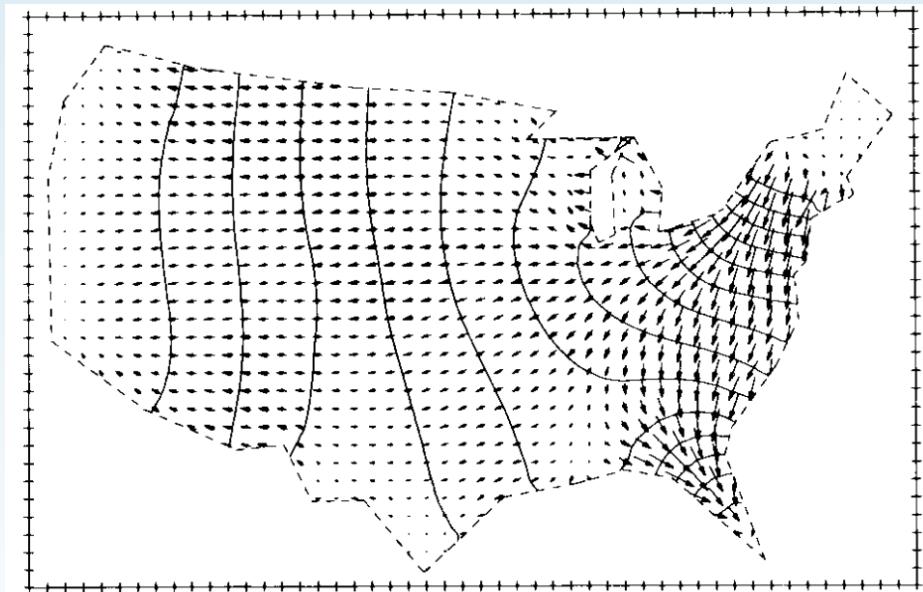
Changes in rates can reflect changes in the numbers of tests performed as well as changes in the prevalence of the virus.



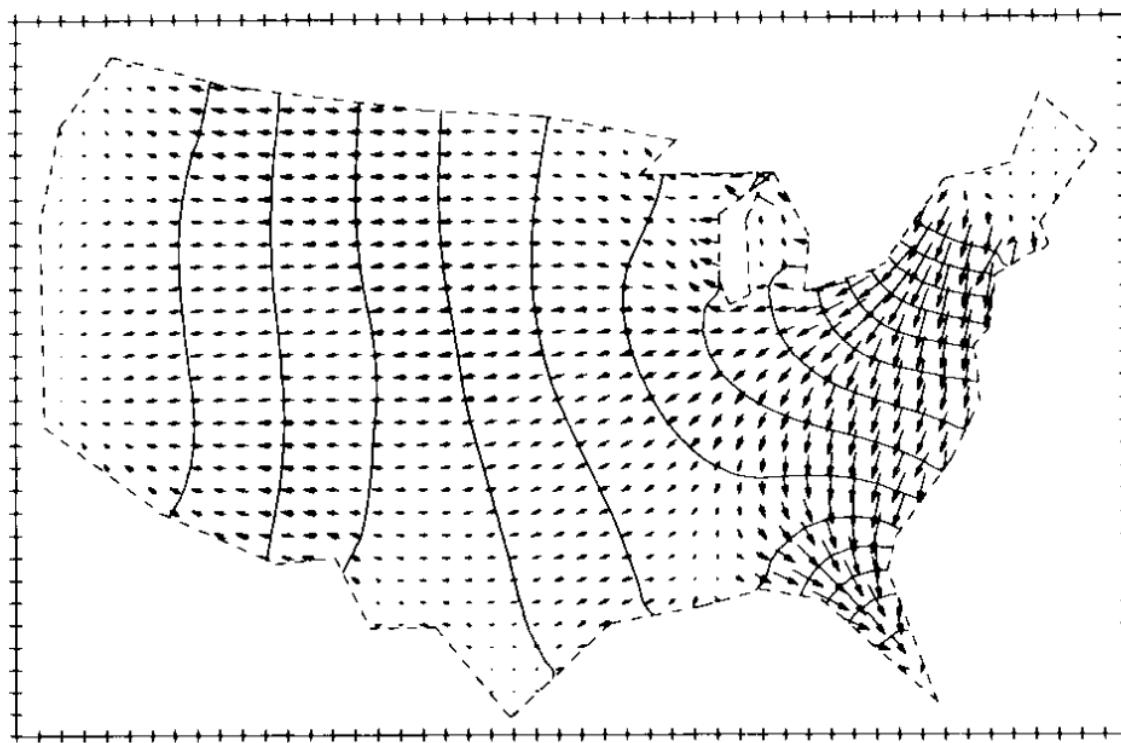
Map by @carlbaker

Data source: Public Health England
Template: tinyurl.com/HCL-hex-cartograms

- Value is shown solely by orientation



Source: Tobler (1981)



MAP 10. Estimated Field of the 1965/1970 Net Population Flow. Obtained as the solution to Poisson's equation. The potential field is indicated by contour lines and the flow is shown by the gradient vectors.

Types of quantitative messages (Few, 2004)

- Nominal comparison
- Time-series
- Ranking
- Part-to-whole
- Deviation
- Frequency distribution
- Correlation

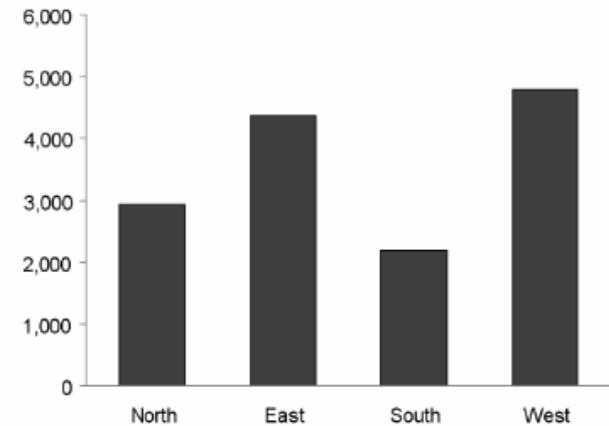
Nominal comparison

Nominal Comparison

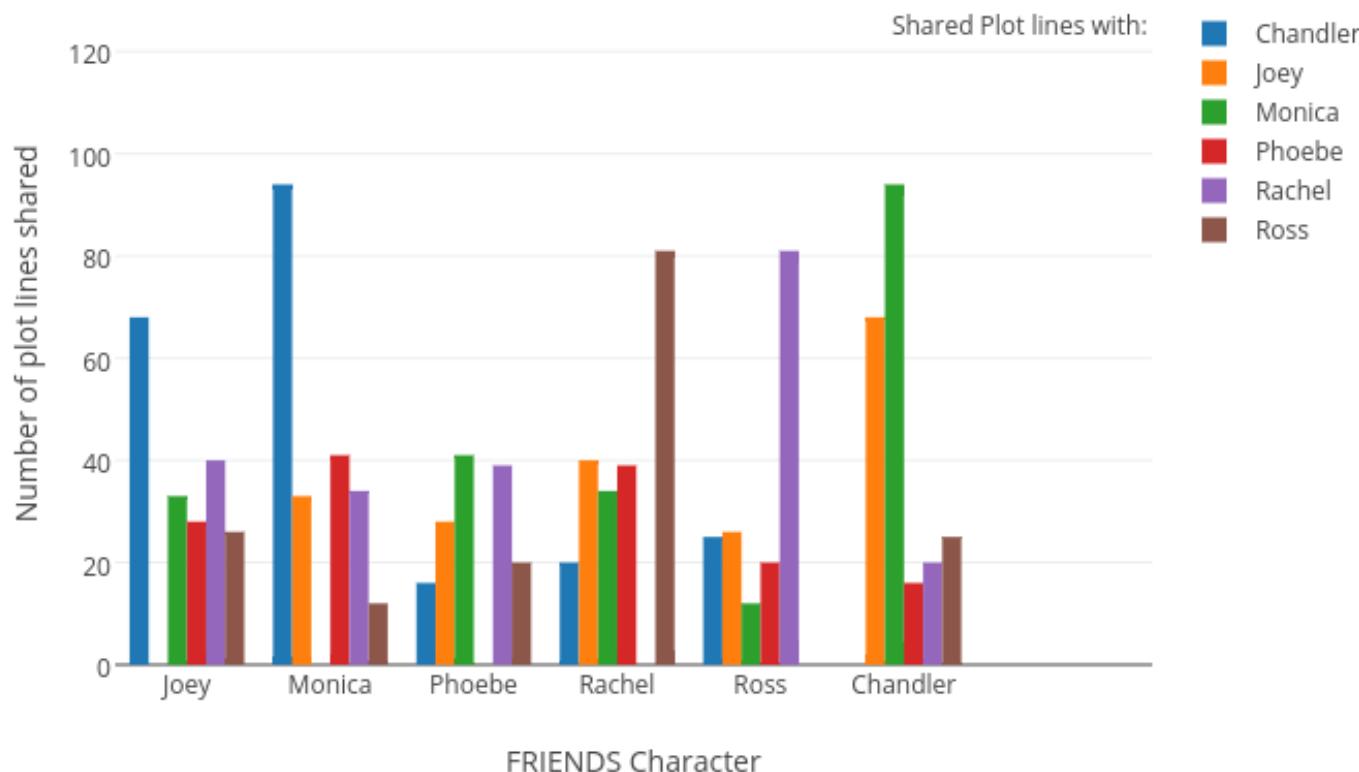
A simple comparison of the categorical subdivisions of one or more measures in no particular order

- Bars only (horizontal or vertical)

Q1 2003 Calls by Region



FRIENDS Plot Line-Interaction Chart



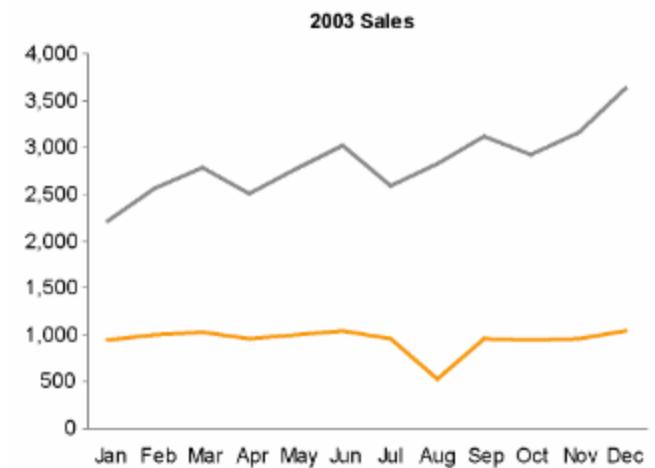
Source: <https://chart-studio.plotly.com/~kirkjtjared/117.embed>

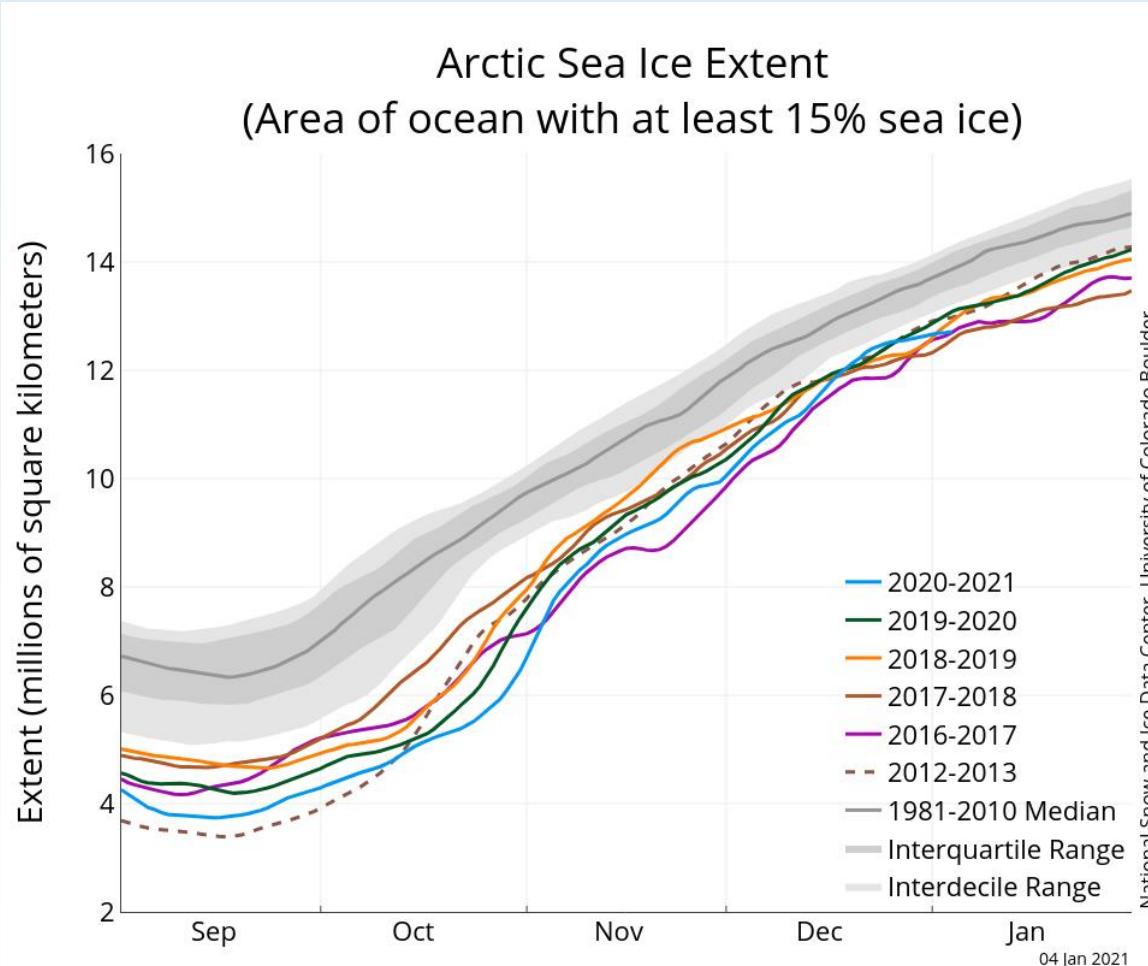
Time series

Time Series

Multiple instances of one or more measures taken at equidistant points in time

- Lines to emphasize overall pattern
- Bars to emphasize individual values
- Points connected by lines to slightly emphasize individual values while still highlighting the overall pattern
- Always place time on the horizontal axis





Source: <http://nsidc.org/arcticseaicenews/>

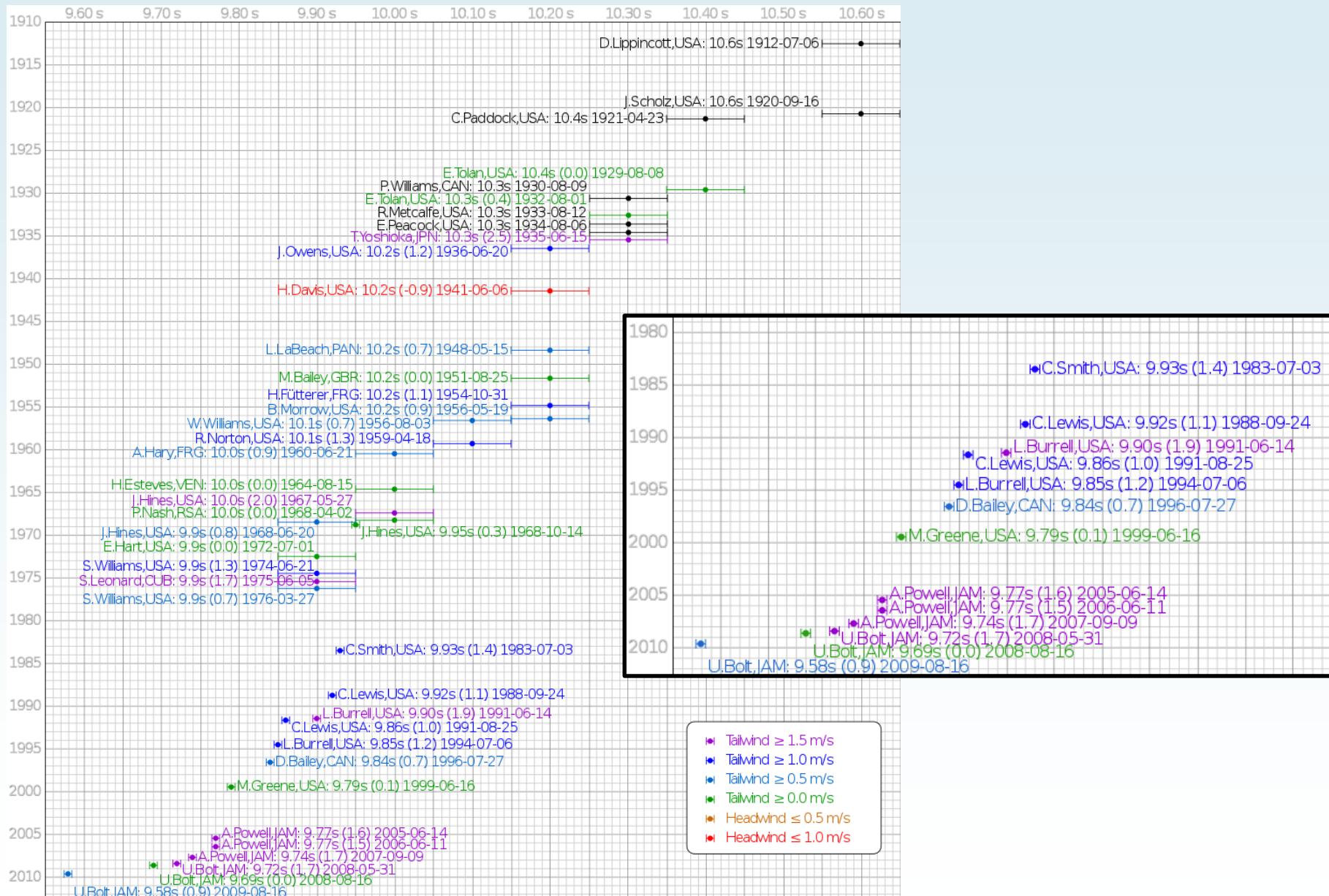
Ranking

Ranking

Categorical subdivisions of a measure ordered by size (either descending or ascending)

- Bars only (horizontal or vertical)
- To highlight high values, sort in descending order
- To highlight low values, sort in ascending order



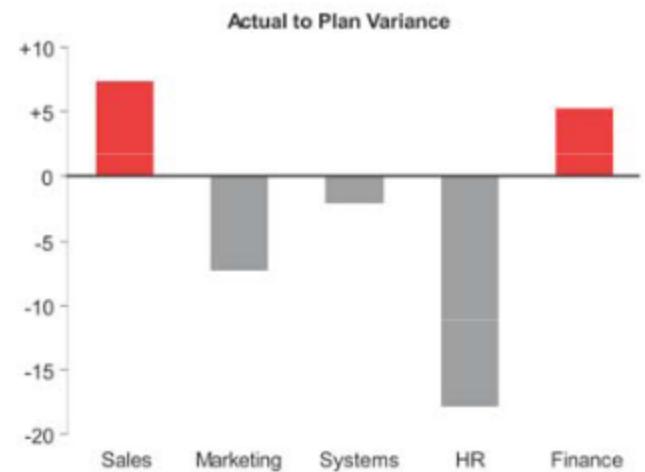


Deviation

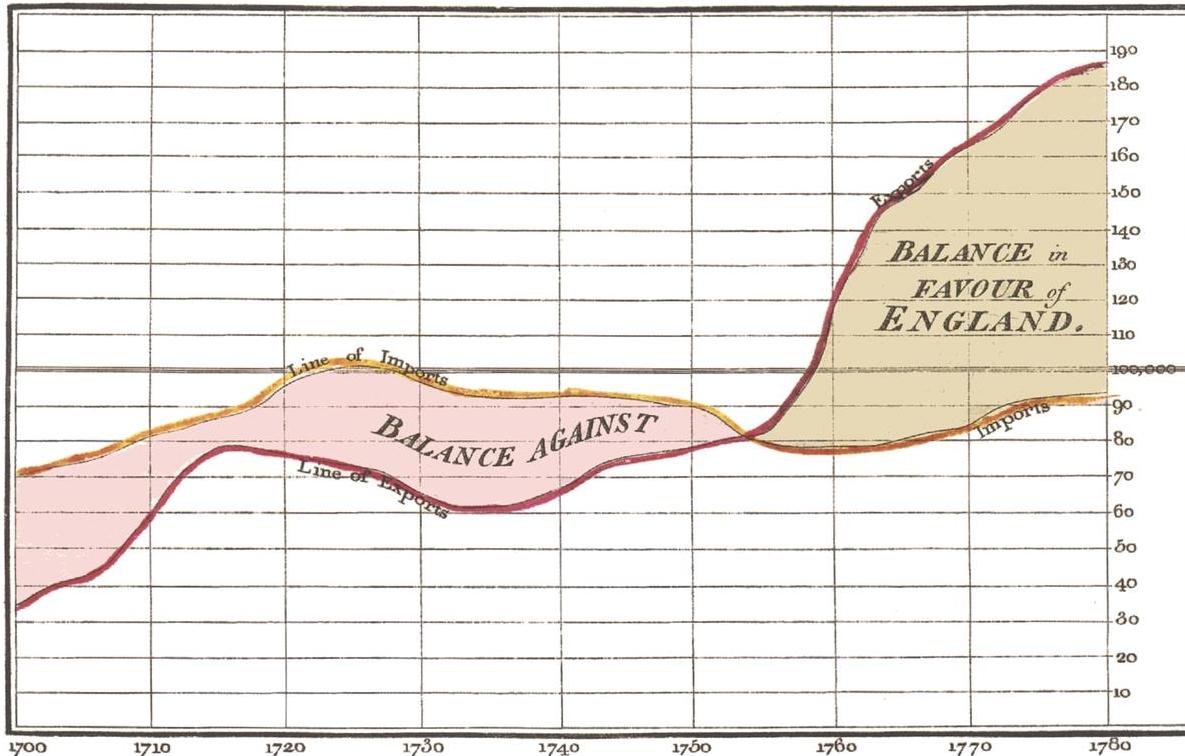
Deviation

Categorical subdivisions of a measure compared to a reference measure, expressed as the differences between them

- Lines to emphasize the overall pattern only when displaying deviation and time-series relationships together
- Points connected by lines to slightly emphasize individual data points while also highlighting the overall pattern when displaying deviation and time-series relationships together
- Bars to emphasize individual values, but limit to vertical bars when a time-series relationship is included
- Always include a reference line to compare the measures of deviation against



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



The Bottom line is divided into Years, the Right hand line into £10,000 each.
 Published as the Act directs, 1st May 1786, by W^m Playfair
 Neale sculpt 352, Strand, London.

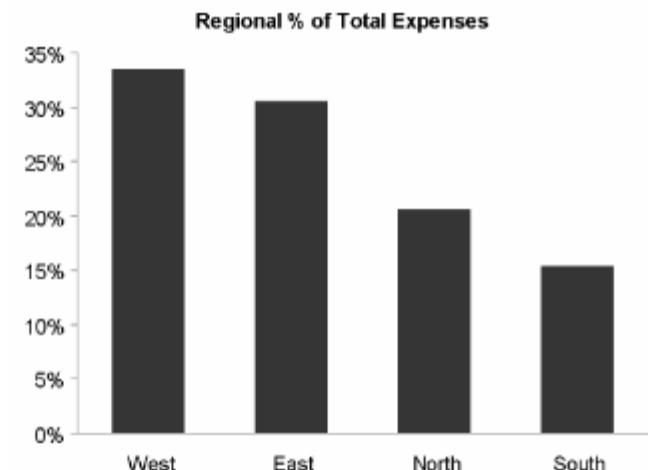
Source (original): Playfair, W (1786)

Part-to-whole

Part-to-Whole

Measures of individual categorical subdivisions as ratios to the whole

- Bars only (horizontal or vertical)
- Use stacked bars only when you must display measures of the whole as well as the parts

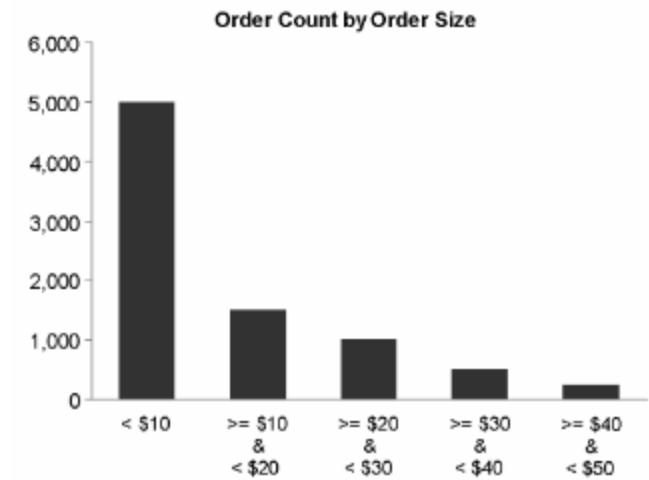


Frequency distribution

Frequency Distribution

Counts of something per categorical subdivisions (intervals) of a quantitative range

- Vertical bars to emphasize individual values (called a *histogram*)
- Lines to emphasize the overall pattern (called a *frequency polygon*)



Variable Description		Representation		
Value	Label		Frequency	% of valid
-9	missing		0	
-8	inapplicable		28,731	
-7	proxy		3,856	
-2	refusal		0	
-1	don't know		3	
1	at least once a week		1,693	13.84%
2	less often than once a week but at least once a month		4,011	32.78%
3	less often than once a month but at least 3 or 4 times a year		3,786	30.94%
4	twice in the last 12 months		1,847	15.09%
5	once in the last 12 months		899	7.35%
Valid				
12236	at least once a week		1,693	13.84%
	less often than once a week but at least once a month		4,011	32.78%
	less often than once a month but at least 3 or 4 times a year		3,786	30.94%
	twice in the last 12 months		1,847	15.09%
	once in the last 12 months		899	7.35%

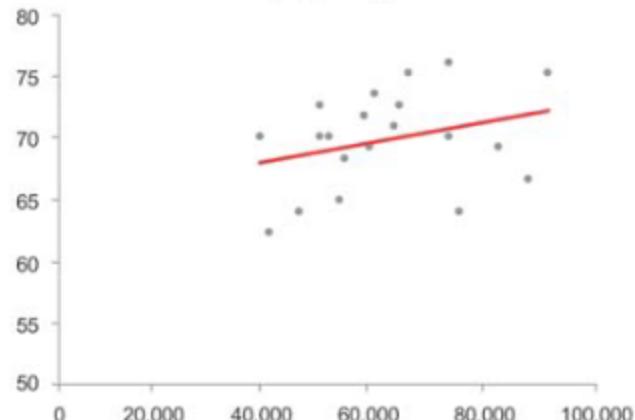
Correlation

Correlation

Comparisons of two paired sets of measures to determine if as one set goes up the other set goes either up or down in a corresponding manner, and if so, how strongly

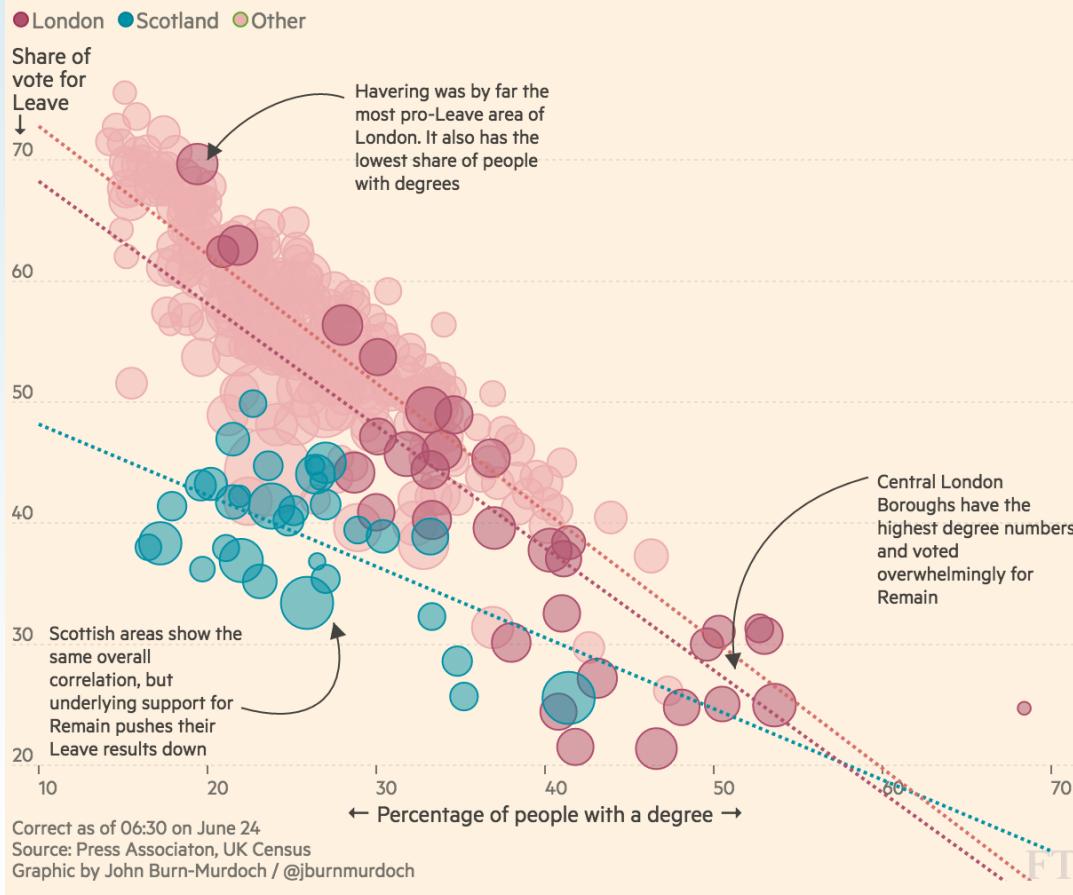
- Points and a trend line in the form of a scatter plot
- Bars may be used, arranged as a *paired bar graph* or a *correlation bar graph*, if scatter plots are unfamiliar
- (Note: For descriptions of these graphs, see my book *Show Me the Numbers*.)

Correlation of Employee Heights and Salaries



A people divided

The strongest correlation between the vote for Leave and any key demographic measure is with the share of people holding a degree. But even here, regional patterns are clear: London Boroughs stand out in the tail on the right, with higher education and low Leave numbers. Scotland follows the overall national trend but is shifted as a whole towards Remain



Source: Financial Times, via:

<https://stats.stackexchange.com/questions/260238/how-should-this-bbc-chart-brexit-correlation-between-education-and-results-hav>

The history of data visualization

- This section identifies some key developments in the history of data visualization
- It is based on:

Friendly M. (2005) Milestones in the History of Data Visualization: A Case Study in Statistical Historiography. In: Weihs C., Gaul W. (eds) Classification — the Ubiquitous Challenge. Studies in Classification, Data Analysis, and Knowledge Organization. Springer, Berlin, Heidelberg.

https://doi.org/10.1007/3-540-28084-7_4

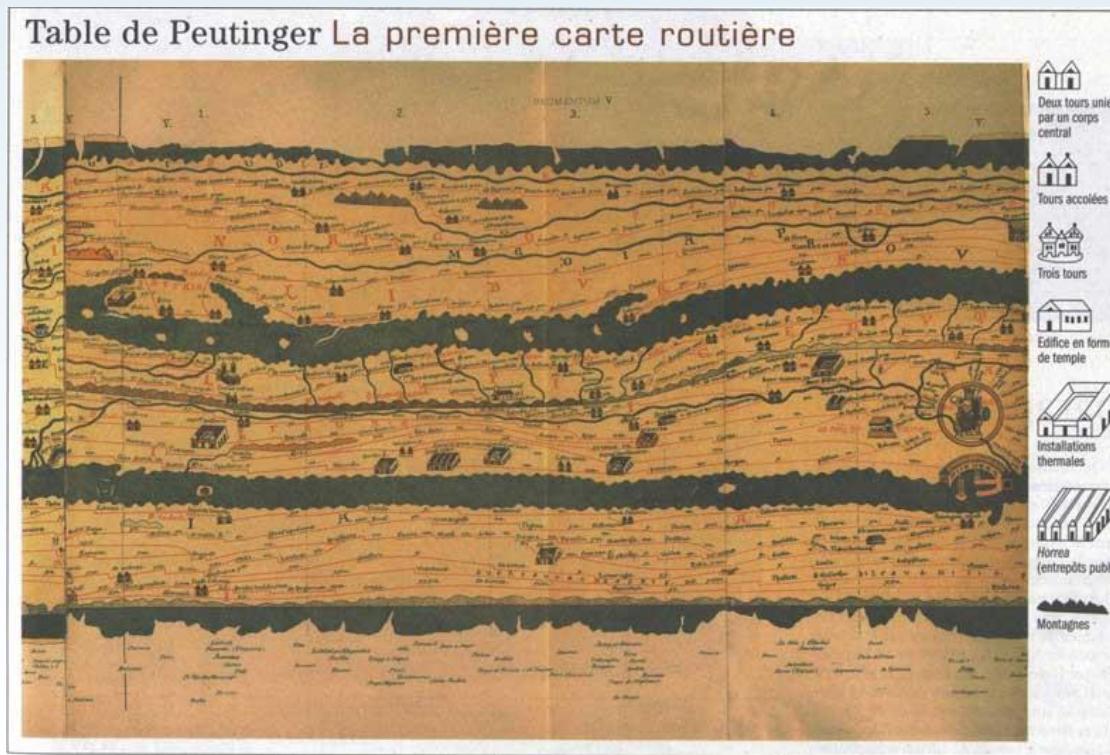
The Milestones website

- <https://www.datavis.ca/milestones/>



Peutinger map, 13th century

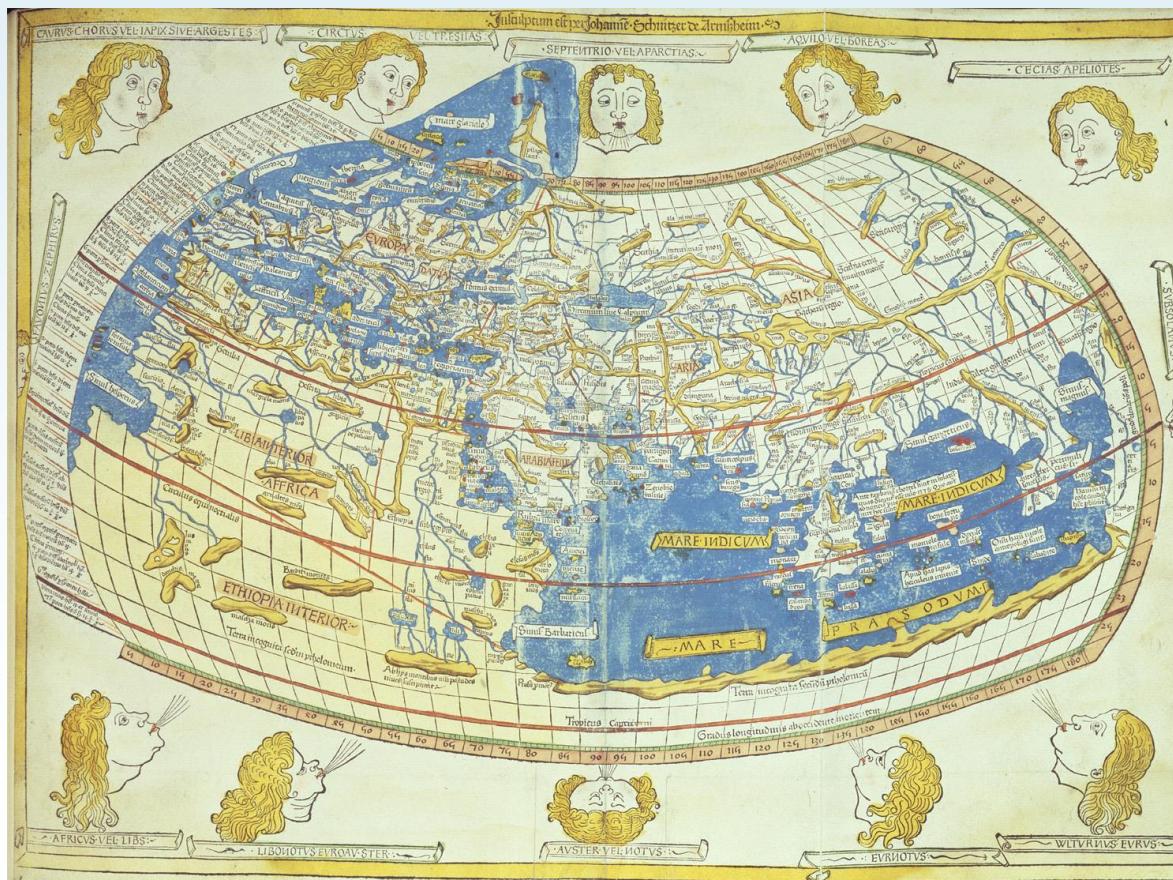
- Based on early maps, possibly dating to c 4th century BC



Source: <https://www.datavis.ca/milestones/index.php?group=pre-1600>

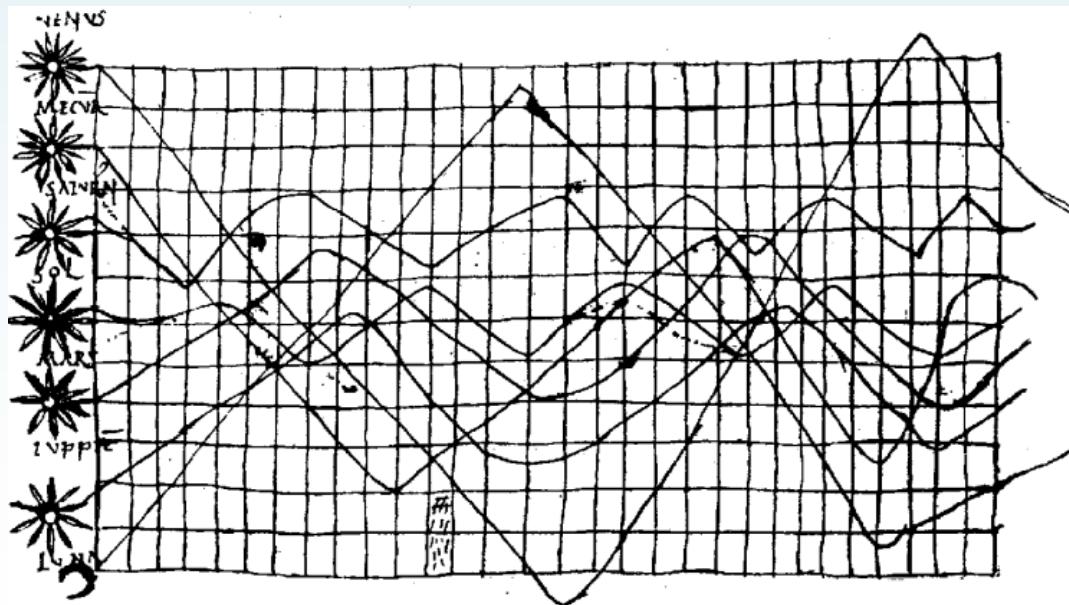
Ptolemy's map, c.150

- Uses latitude and longitude



Unknown, c. 950

- "Earliest known attempt to show changing values graphically (positions of the sun, moon, and planets throughout the year)"



Source: <https://www.datavis.ca/milestones/index.php?group=pre-1600#>

Langren, 1644

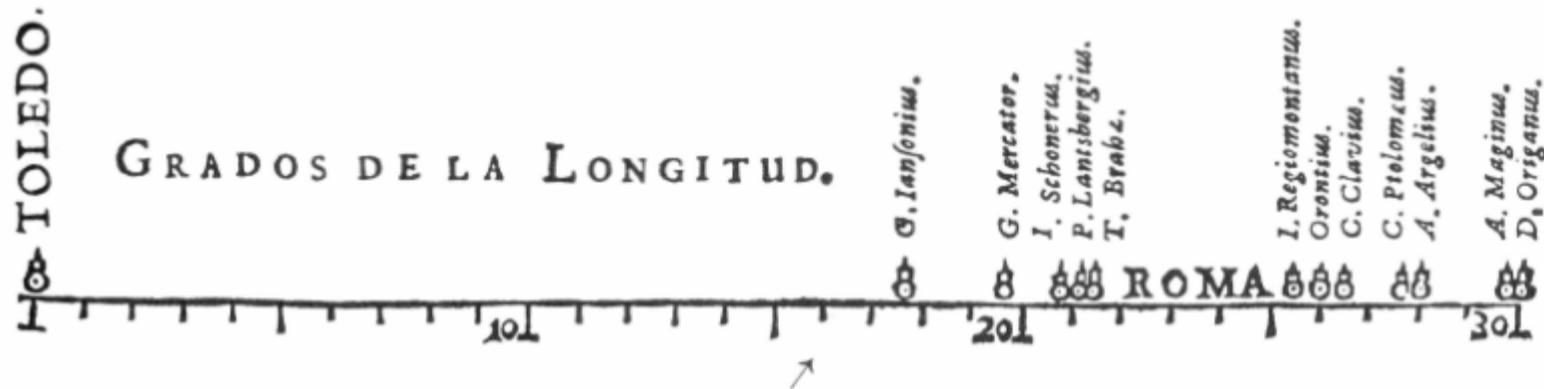


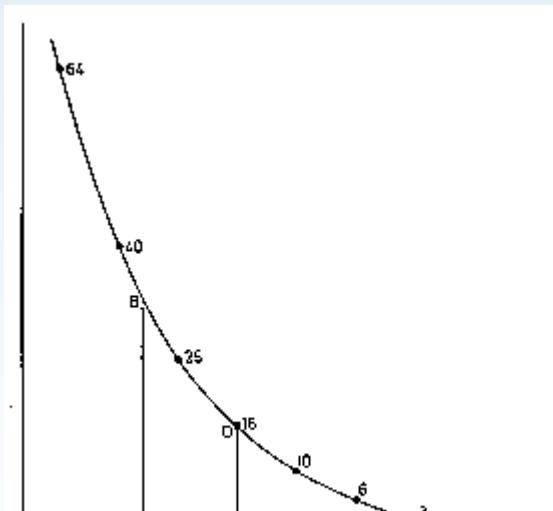
Fig. 1. Langren's 1644 graph of determinations of the distance, in longitude, from Toledo to Rome. The correct distance is $16^{\circ}30'$. Source: Tufte (1997, p. 15)

Life tables – Graunt, 1662

Graphical representation – Huygens, 1669

Buried of all Diseases in the Year 1592.	Total / Plague	Buried of all Diseases in the Year 1603.	Total / Plague	Buried of all Diseases in the Year 1625.	Total / Plague	Buried of all Diseases in the Year 1630.	Total / Plague	Buried of all Diseases in the Year 1636.	Total / Plague
March 17	230 / 3	March 17	108 / 3	March 17	261 / 5	June 24	205 / 19	April 7	196 / 4
March 18	230 / 3	March 18	108 / 3	March 18	261 / 5	June 25	205 / 19	April 8	196 / 4
March 21	210 / 29	March 21	88 / 29	March 21	243 / 1	July 8	217 / 43	April 9	205 / 4
March 24	210 / 29	March 24	88 / 29	March 24	243 / 1	July 9	217 / 43	April 10	205 / 4
April 7	305 / 27	April 7	60 / 27	April 7	239 / 4	July 15	250 / 50	April 11	305 / 4
April 8	305 / 27	April 8	60 / 27	April 8	239 / 4	July 16	250 / 50	April 12	305 / 4
April 11	295 / 27	April 11	52 / 27	April 11	239 / 4	July 17	250 / 50	April 13	295 / 4
April 12	295 / 27	April 12	52 / 27	April 12	239 / 4	July 18	250 / 50	April 14	295 / 4
April 15	295 / 27	April 15	52 / 27	April 15	239 / 4	July 19	250 / 50	April 15	295 / 4
April 16	295 / 27	April 16	52 / 27	April 16	239 / 4	July 20	250 / 50	April 16	295 / 4
May 1	350 / 28	May 1	90 / 28	May 1	285 / 10	August 5	242 / 65	April 17	350 / 4
May 2	350 / 28	May 2	90 / 28	May 2	285 / 10	August 6	242 / 65	April 18	350 / 4
May 3	350 / 28	May 3	90 / 28	May 3	285 / 10	August 7	242 / 65	April 19	350 / 4
May 4	350 / 28	May 4	90 / 28	May 4	285 / 10	August 8	242 / 65	April 20	350 / 4
May 5	350 / 28	May 5	90 / 28	May 5	285 / 10	August 9	242 / 65	April 21	350 / 4
May 6	350 / 28	May 6	90 / 28	May 6	285 / 10	August 10	242 / 65	April 22	350 / 4
May 7	350 / 28	May 7	90 / 28	May 7	285 / 10	August 11	242 / 65	April 23	350 / 4
May 8	350 / 28	May 8	90 / 28	May 8	285 / 10	August 12	242 / 65	April 24	350 / 4
May 9	350 / 28	May 9	90 / 28	May 9	285 / 10	August 13	242 / 65	April 25	350 / 4
May 10	350 / 28	May 10	90 / 28	May 10	285 / 10	August 14	242 / 65	April 26	350 / 4
May 11	350 / 28	May 11	90 / 28	May 11	285 / 10	August 15	242 / 65	April 27	350 / 4
May 12	350 / 28	May 12	90 / 28	May 12	285 / 10	August 16	242 / 65	April 28	350 / 4
May 13	350 / 28	May 13	90 / 28	May 13	285 / 10	August 17	242 / 65	April 29	350 / 4
May 14	350 / 28	May 14	90 / 28	May 14	285 / 10	August 18	242 / 65	April 30	350 / 4
May 15	350 / 28	May 15	90 / 28	May 15	285 / 10	August 19	242 / 65	May 1	350 / 4
May 16	350 / 28	May 16	90 / 28	May 16	285 / 10	August 20	242 / 65	May 2	350 / 4
May 17	350 / 28	May 17	90 / 28	May 17	285 / 10	August 21	242 / 65	May 3	350 / 4
May 18	350 / 28	May 18	90 / 28	May 18	285 / 10	August 22	242 / 65	May 4	350 / 4
May 19	350 / 28	May 19	90 / 28	May 19	285 / 10	August 23	242 / 65	May 5	350 / 4
May 20	350 / 28	May 20	90 / 28	May 20	285 / 10	August 24	242 / 65	May 6	350 / 4
May 21	350 / 28	May 21	90 / 28	May 21	285 / 10	August 25	242 / 65	May 7	350 / 4
May 22	350 / 28	May 22	90 / 28	May 22	285 / 10	August 26	242 / 65	May 8	350 / 4
May 23	350 / 28	May 23	90 / 28	May 23	285 / 10	August 27	242 / 65	May 9	350 / 4
May 24	350 / 28	May 24	90 / 28	May 24	285 / 10	August 28	242 / 65	May 10	350 / 4
May 25	350 / 28	May 25	90 / 28	May 25	285 / 10	August 29	242 / 65	May 11	350 / 4
May 26	350 / 28	May 26	90 / 28	May 26	285 / 10	August 30	242 / 65	May 12	350 / 4
May 27	350 / 28	May 27	90 / 28	May 27	285 / 10	August 31	242 / 65	May 13	350 / 4
May 28	350 / 28	May 28	90 / 28	May 28	285 / 10	September 1	242 / 65	May 14	350 / 4
May 29	350 / 28	May 29	90 / 28	May 29	285 / 10	September 2	242 / 65	May 15	350 / 4
May 30	350 / 28	May 30	90 / 28	May 30	285 / 10	September 3	242 / 65	May 16	350 / 4
May 31	350 / 28	May 31	90 / 28	May 31	285 / 10	September 4	242 / 65	May 17	350 / 4
June 1	350 / 28	June 1	90 / 28	June 1	285 / 10	September 5	242 / 65	May 18	350 / 4
June 2	350 / 28	June 2	90 / 28	June 2	285 / 10	September 6	242 / 65	May 19	350 / 4
June 3	350 / 28	June 3	90 / 28	June 3	285 / 10	September 7	242 / 65	May 20	350 / 4
June 4	350 / 28	June 4	90 / 28	June 4	285 / 10	September 8	242 / 65	May 21	350 / 4
June 5	350 / 28	June 5	90 / 28	June 5	285 / 10	September 9	242 / 65	May 22	350 / 4
June 6	350 / 28	June 6	90 / 28	June 6	285 / 10	September 10	242 / 65	May 23	350 / 4
June 7	350 / 28	June 7	90 / 28	June 7	285 / 10	September 11	242 / 65	May 24	350 / 4
June 8	350 / 28	June 8	90 / 28	June 8	285 / 10	September 12	242 / 65	May 25	350 / 4
June 9	350 / 28	June 9	90 / 28	June 9	285 / 10	September 13	242 / 65	May 26	350 / 4
June 10	350 / 28	June 10	90 / 28	June 10	285 / 10	September 14	242 / 65	May 27	350 / 4
June 11	350 / 28	June 11	90 / 28	June 11	285 / 10	September 15	242 / 65	May 28	350 / 4
June 12	350 / 28	June 12	90 / 28	June 12	285 / 10	September 16	242 / 65	May 29	350 / 4
June 13	350 / 28	June 13	90 / 28	June 13	285 / 10	September 17	242 / 65	May 30	350 / 4
June 14	350 / 28	June 14	90 / 28	June 14	285 / 10	September 18	242 / 65	May 31	350 / 4
June 15	350 / 28	June 15	90 / 28	June 15	285 / 10	September 19	242 / 65	June 1	350 / 4
June 16	350 / 28	June 16	90 / 28	June 16	285 / 10	September 20	242 / 65	June 2	350 / 4
June 17	350 / 28	June 17	90 / 28	June 17	285 / 10	September 21	242 / 65	June 3	350 / 4
June 18	350 / 28	June 18	90 / 28	June 18	285 / 10	September 22	242 / 65	June 4	350 / 4
June 19	350 / 28	June 19	90 / 28	June 19	285 / 10	September 23	242 / 65	June 5	350 / 4
June 20	350 / 28	June 20	90 / 28	June 20	285 / 10	September 24	242 / 65	June 6	350 / 4
June 21	350 / 28	June 21	90 / 28	June 21	285 / 10	September 25	242 / 65	June 7	350 / 4
June 22	350 / 28	June 22	90 / 28	June 22	285 / 10	September 26	242 / 65	June 8	350 / 4
June 23	350 / 28	June 23	90 / 28	June 23	285 / 10	September 27	242 / 65	June 9	350 / 4
June 24	350 / 28	June 24	90 / 28	June 24	285 / 10	September 28	242 / 65	June 10	350 / 4
June 25	350 / 28	June 25	90 / 28	June 25	285 / 10	September 29	242 / 65	June 11	350 / 4
June 26	350 / 28	June 26	90 / 28	June 26	285 / 10	September 30	242 / 65	June 12	350 / 4
June 27	350 / 28	June 27	90 / 28	June 27	285 / 10	September 31	242 / 65	June 13	350 / 4
June 28	350 / 28	June 28	90 / 28	June 28	285 / 10	October 1	242 / 65	June 14	350 / 4
June 29	350 / 28	June 29	90 / 28	June 29	285 / 10	October 2	242 / 65	June 15	350 / 4
June 30	350 / 28	June 30	90 / 28	June 30	285 / 10	October 3	242 / 65	June 16	350 / 4
June 31	350 / 28	June 31	90 / 28	June 31	285 / 10	October 4	242 / 65	June 17	350 / 4
The Total of all that have been buried is	35086	The Total of all that have been buried is	35086	The Total of all that have been buried is	35086	The Total of all that have been buried is	35086	The Total of all that have been buried is	35086
Whereof the Plague	35086	Whereof the Plague	35086	Whereof the Plague	35086	Whereof the Plague	35086	Whereof the Plague	35086
The Total of all the Burials this year is	35086	The Total of all the Burials this year is	35086	The Total of all the Burials this year is	35086	The Total of all the Burials this year is	35086	The Total of all the Burials this year is	35086
Whereof of the Pt	35086	Whereof of the Pt	35086	Whereof of the Pt	35086	Whereof of the Pt	35086	Whereof of the Pt	35086
The Total of all the Burials this year is	35086	The Total of all the Burials this year is	35086	The Total of all the Burials this year is	35086	The Total of all the Burials this year is	35086	The Total of all the Burials this year is	35086
Whereof of the Pt	35086	Whereof of the Pt	35086	Whereof of the Pt	35086	Whereof of the Pt	35086	Whereof of the Pt	35086

Place this Table at page 406.



Source: <https://www.datavis.ca/milestones/index.php?group=1600s>

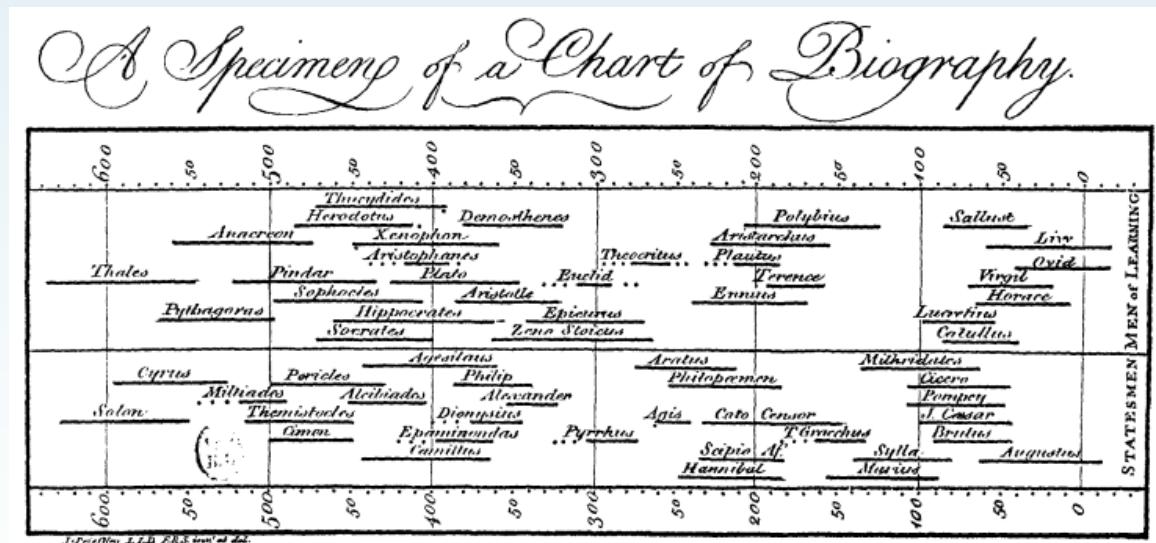
Halley, 1701



- "Contour maps showing curves of equal value (an isogonic map, lines of equal magnetic declination for the world, possibly the first contour map of a data-based variable)"

Priestly, 1765 Historical timeline

- "Historical timeline (life spans of 2,000 famous people, 1200 B.C. to 1750 A.D.), quantitative comparison by means of bars"



William Playfair

"William Playfair (1759–1823) is widely considered the inventor of most of the graphical forms widely used today— first the line graph and bar chart (Playfair, 1786), later the pie chart and circle graph (Playfair, 1801). A somewhat later graph(Playfair, 1821), shown in Figure 2, exemplifies the best that Playfair had to offer with these graphic forms. Playfair used three parallel timeseries to show the price of wheat, weekly wages, and reigning monarch over a~250 year span from 1565to 1820, and used this graph to argue that workers had become better off in the most recent years"

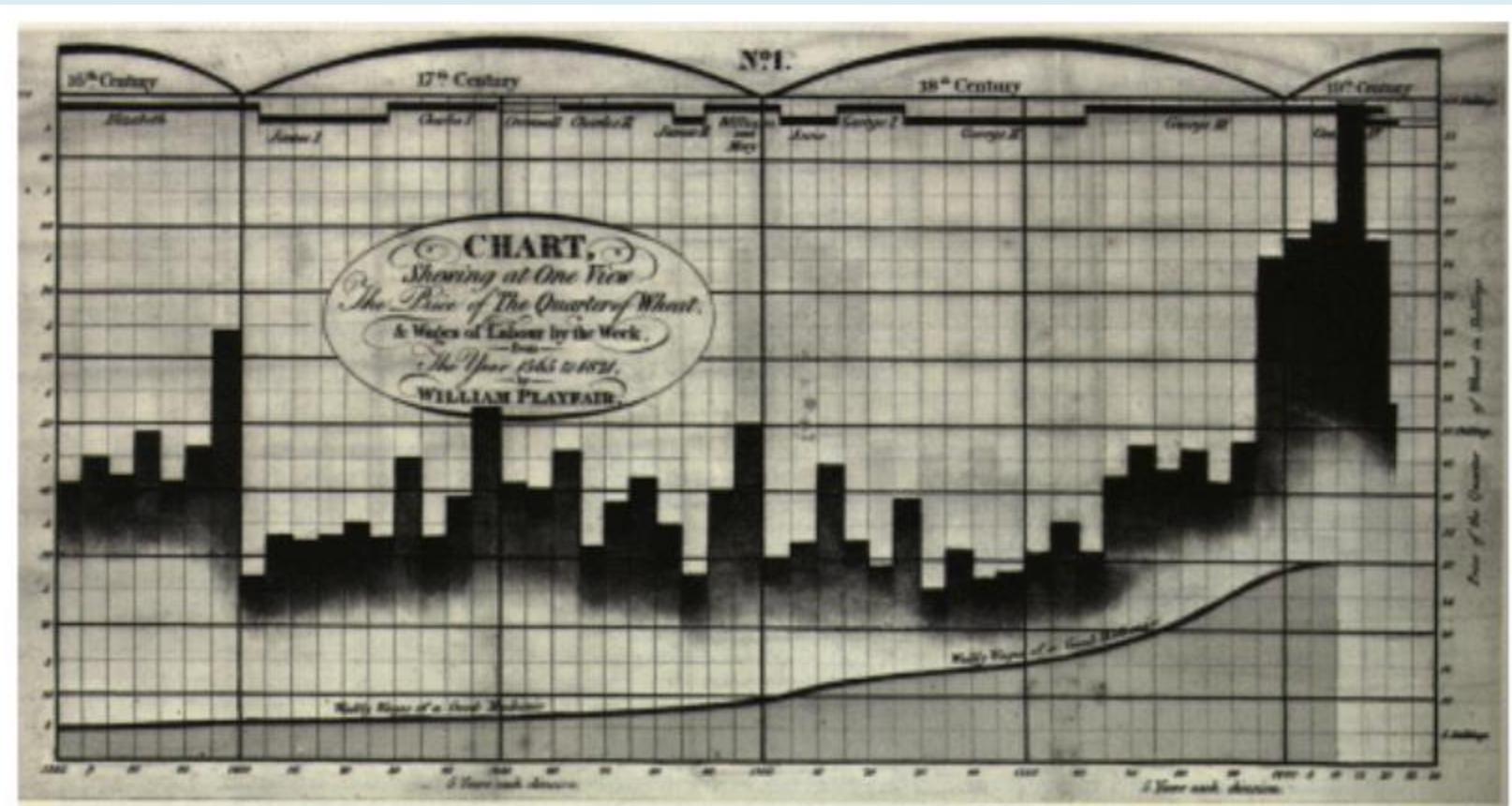
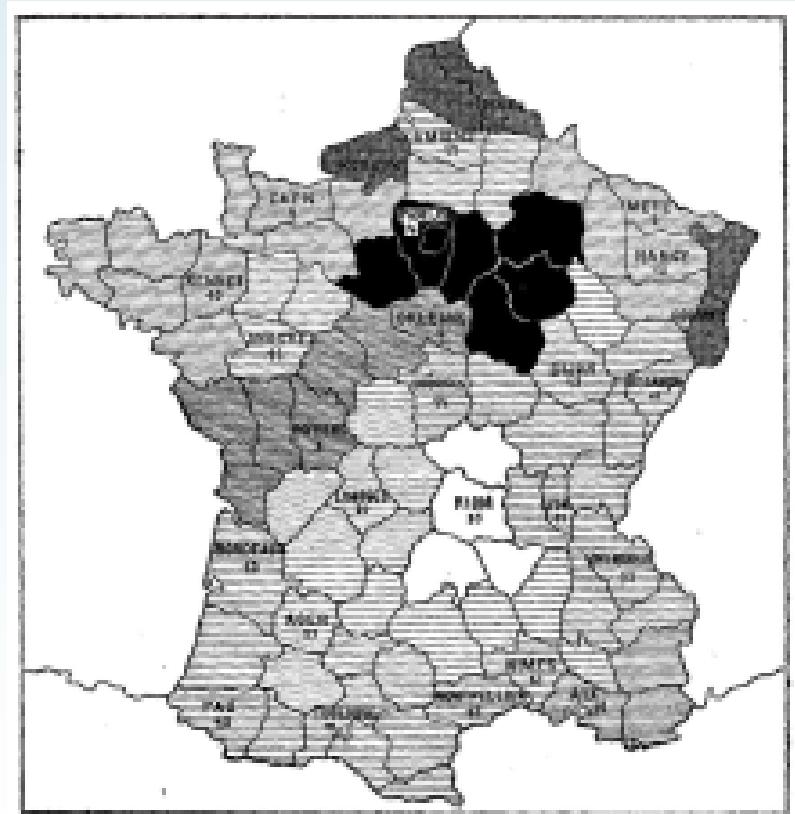


Fig. 2. William Playfair's 1821 time series graph of prices, wages, and ruling monarch over a 250 year period. *Source:* Playfair (1821), image from Tufte (1983, p. 34)

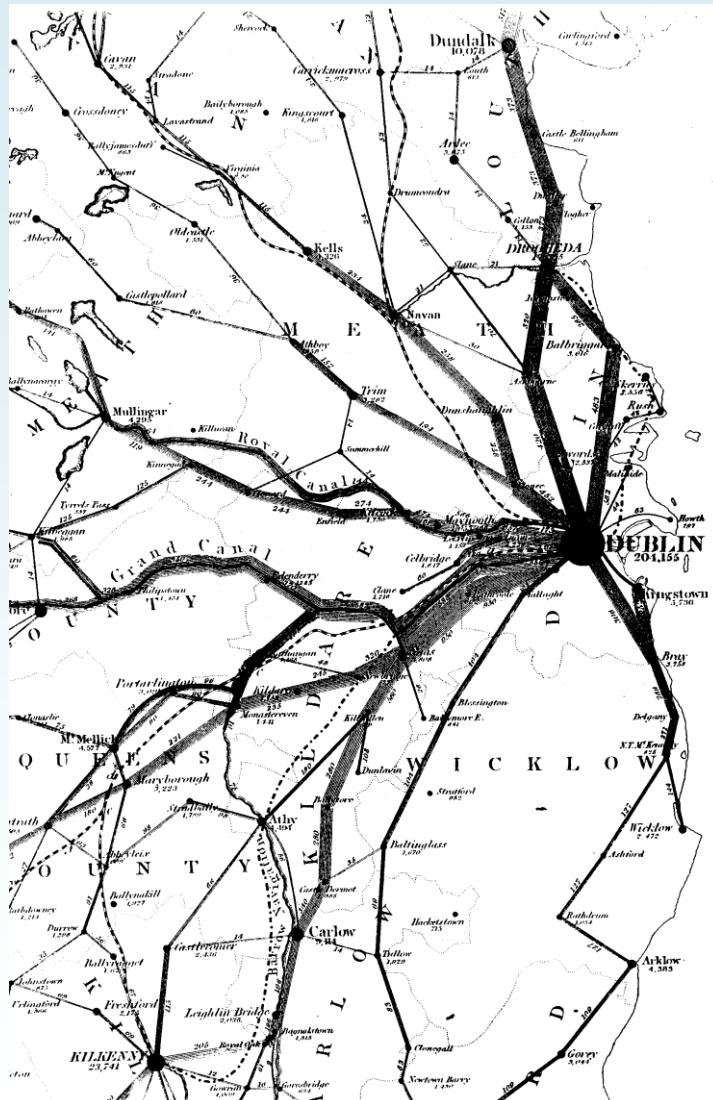
Dupin, 1819

- First choropleth map?



Harness, 1837

- First flow map



Minard, 1844

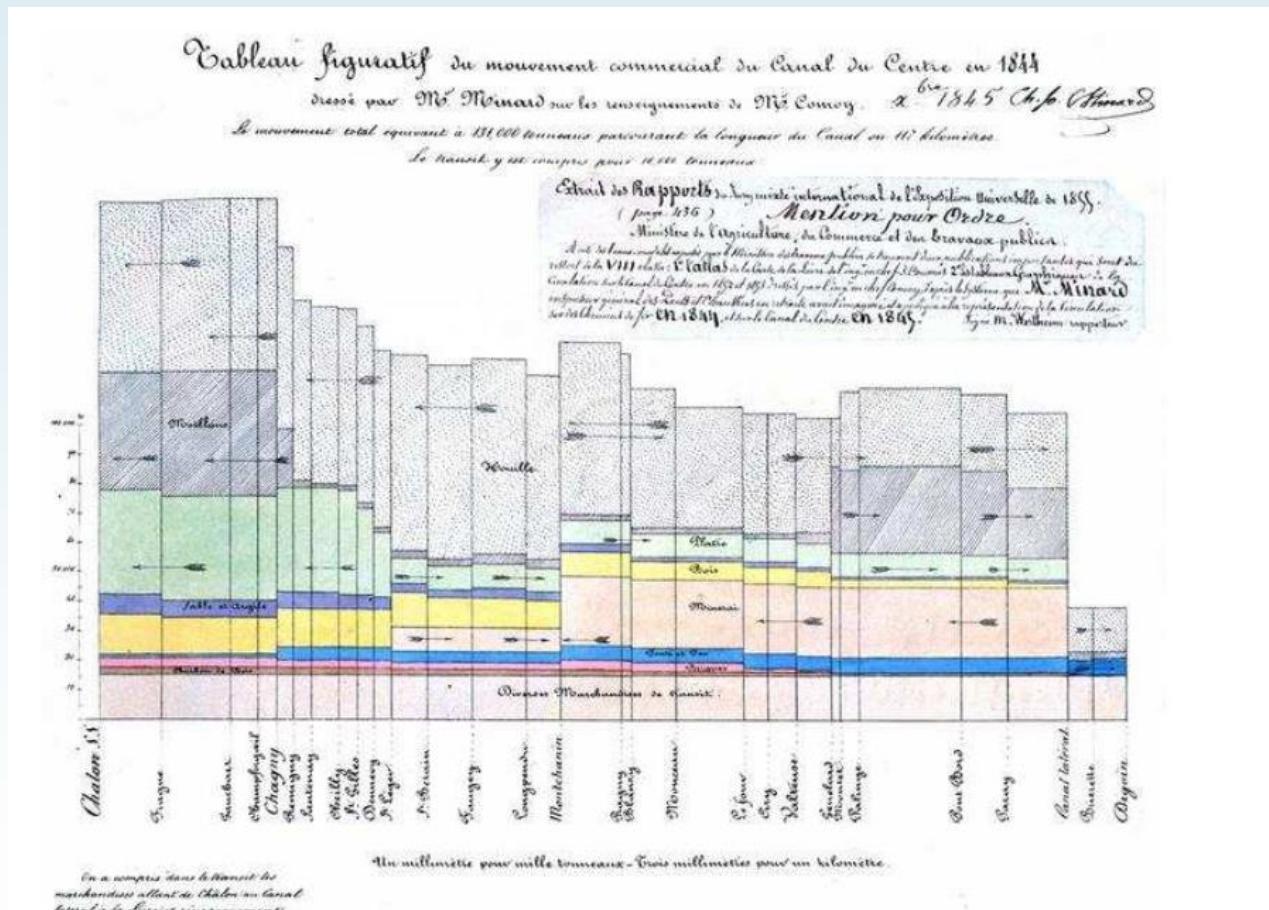


Fig. 3. Minard's Tableau Graphique, showing the transportation of commercial goods along the Canal du Centre (Chalon–Dijon). Intermediate stops are spaced by distance, and each bar is divided by type of goods, so the area of each tile represents the cost of transport. Arrows show the direction of transport. *Source:* ENPC:5860/C351 (Col. et cliché ENPC; used by permission)

Friendly

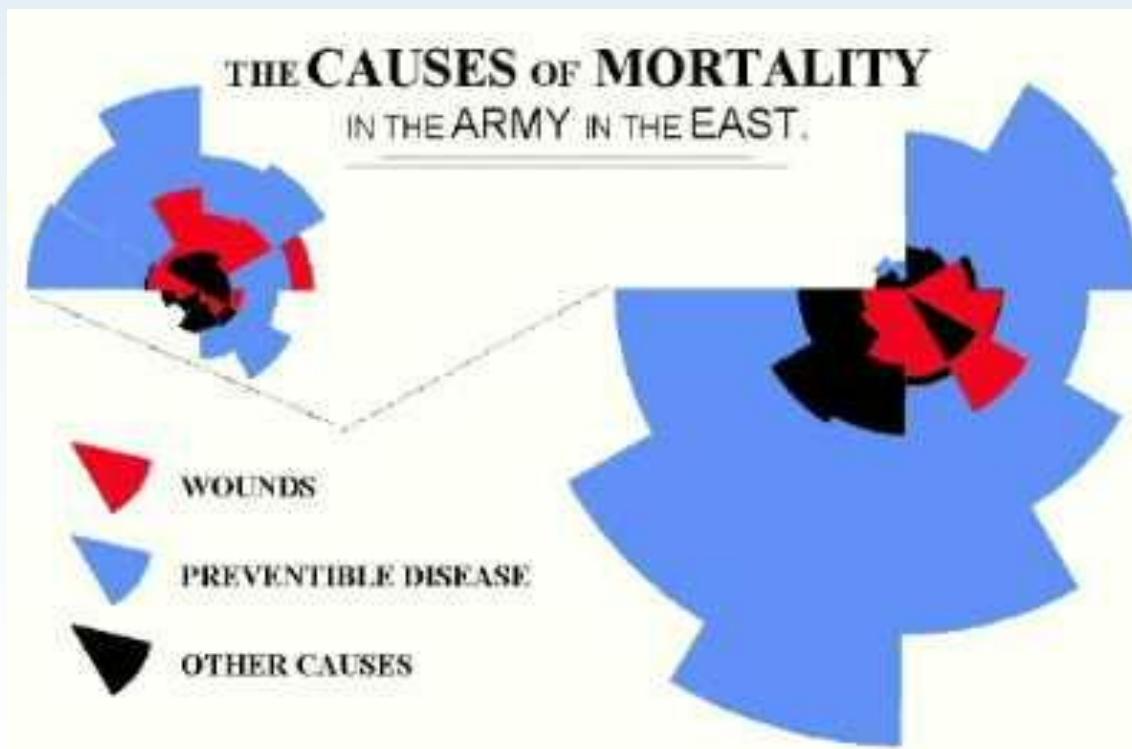
Period	Description
1600-1699	Measurement and theory
1700-1799	New graphic forms
1800-1850	Beginnings of modern graphics
1850-1899	The Golden Age of statistical graphics
1900-1950	The modern dark ages
1950-1975	Re-birth of data visualization

Snow, 1855 Cholera map



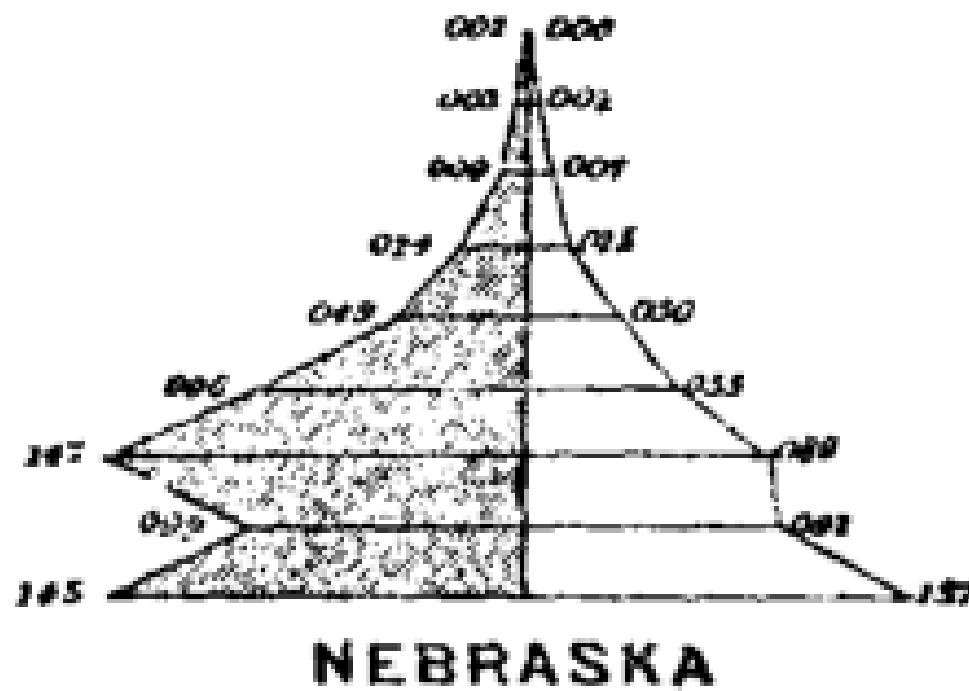
Nightingale, 1857

Causes of mortality

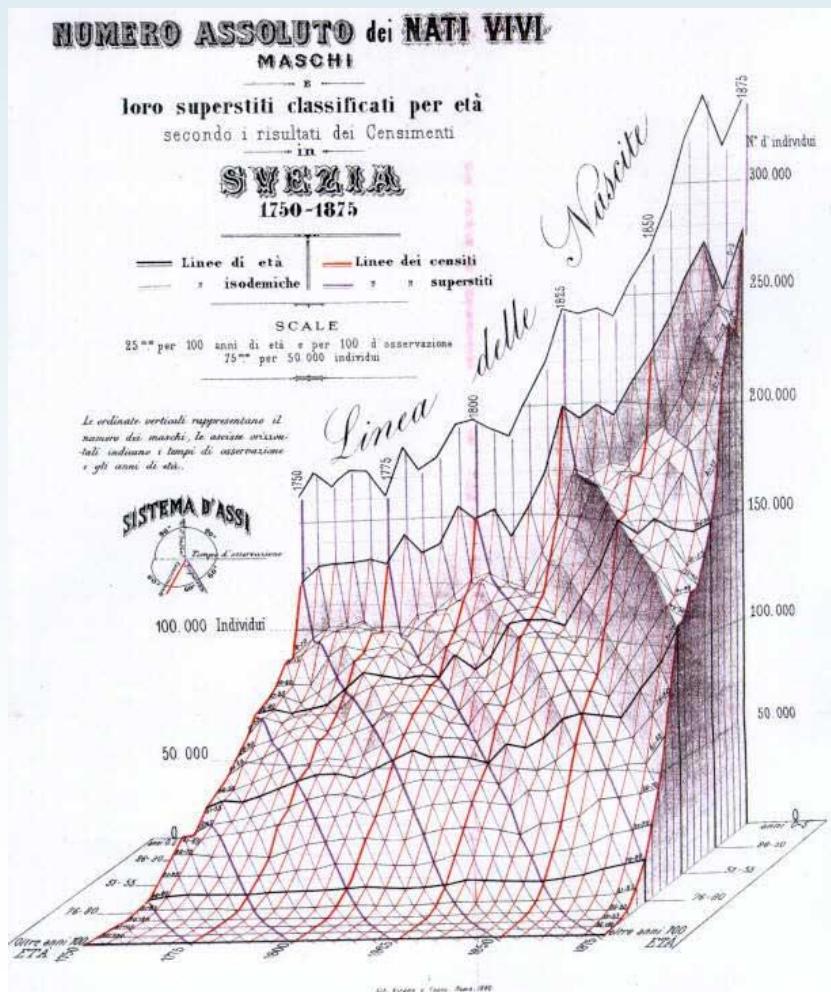


Walker, 1874

Population pyramid

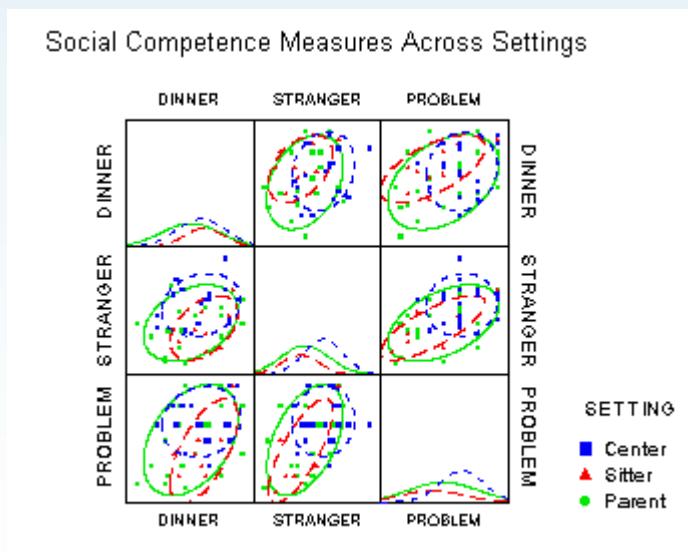


Perozzo, 1879



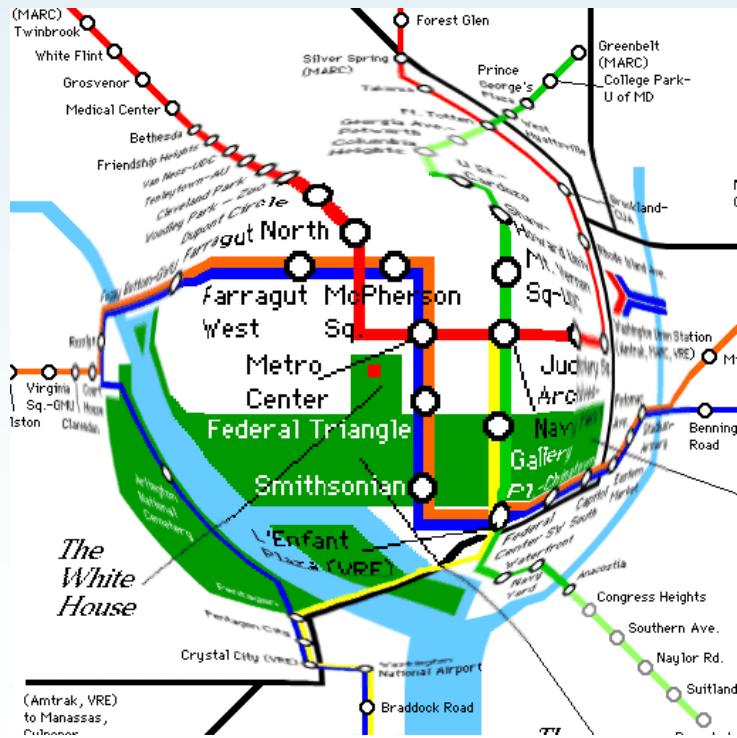
"Stereogram (three-dimensional population pyramid) modeled on actual data (Swedish census, 1750--1875)"

Hartigan, 1975



"Scatterplot matrix, the idea of plotting all pairwise scatterplots for n variables in a tabular display"

Furnas, 1981



"Fisheye view: an idea to provide focus and greater detail in areas of interest of a large amount of information, while retaining the surrounding context in much less detail"

Rosling, 2006

