

9 Visualization Toolbox

There are many types of data visualizations and many variations on each type. These data visualizations can collectively be thought of as tools in the data visualizer’s tool box, and good data visualizers will be as familiar with them as a master wood worker is with their tools of their trade.

The only way this can come about is through practicing with each type and variation of visualization, using either teaching datasets or by finding real world use cases.

It also helps to have a good organization to the tool box, which can serve to direct the visualizer to the right general type of visualization.

As with organizing a tool box, there is no one right way to group data visualizations together, and to some extent it is a matter of personal preference. However, some experts in this field have made efforts on this front, and it’s possible to see some commonalities.

As seen in Figure 9.1, one approach to organizing data visualizations is to consider which ones best highlight:

- a **relationship** – show a connection or correlation between two or more variables,¹ such as the impact of an aging population on health care;
- a **comparison** – set some variables apart from others, and display how those two variables interact, such as the number of fans attending hockey games for different teams in a season;
- a **composition** – collect different types of information that make up a whole and display them together, such as the various search terms that visitors used to land on your site, or how many visitors came from various sources (links, search engines, or direct traffic), and
- a **distribution** – lay out a collection of related or unrelated information to see how it correlates (if at all), and to understand if there’s any interaction between the variables, such as the number of bugs reported during each month after a new software release.

However, this is not the only way to think about data visualizations. Some practitioners have broken down these categories further, as shown in Figure 9.2. As yet another alternative, it’s possible to consider which combinations of the 5W questions (**who, what, when, where, and how/why**) certain data visualizations are best suited to display, as was shown in Figure 2.14.

9.1 Exploration Visualizations	179
Scatterplots	179
Rug Charts and Histograms	180
Two-Way Tables	182
9.2 Presentation Visualizations	183
Text Blocks	183
Tables	184
Line Graphs	184
Bar Charts	185
9.3 The Rest of The Landscape	188
Maps and Heat Maps	188
Bubble Charts	189
Small Multiples	189
Area Charts and Treemaps .	191
Text Visualizations	191
Parallel Coordinates	193
Trees and Networks	193
Animated Visualizations .	193
9.4 Misc. & Charts to Avoid . .	197
Chernoff Faces	197
Alluvial Diagrams	197
Charts to Avoid	198

1: Also called dimensions, axes, factors, etc.

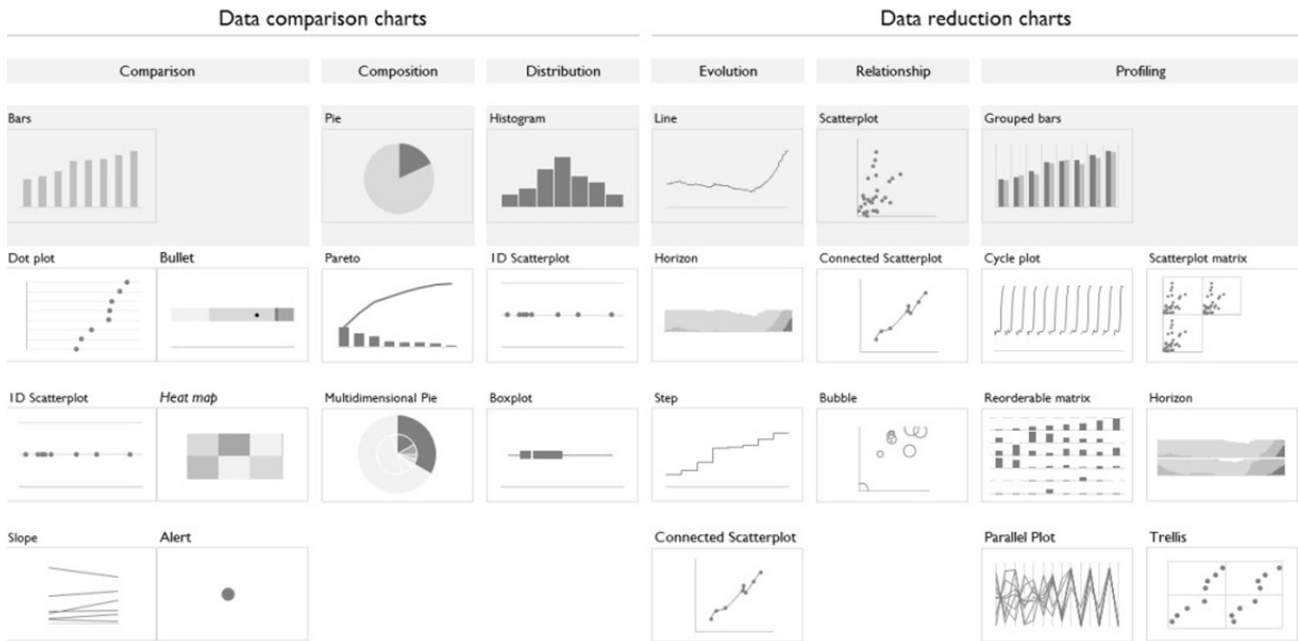


Figure 9.1: A Classification of chart types, based on visualization objectives [J. Camoes et al.].

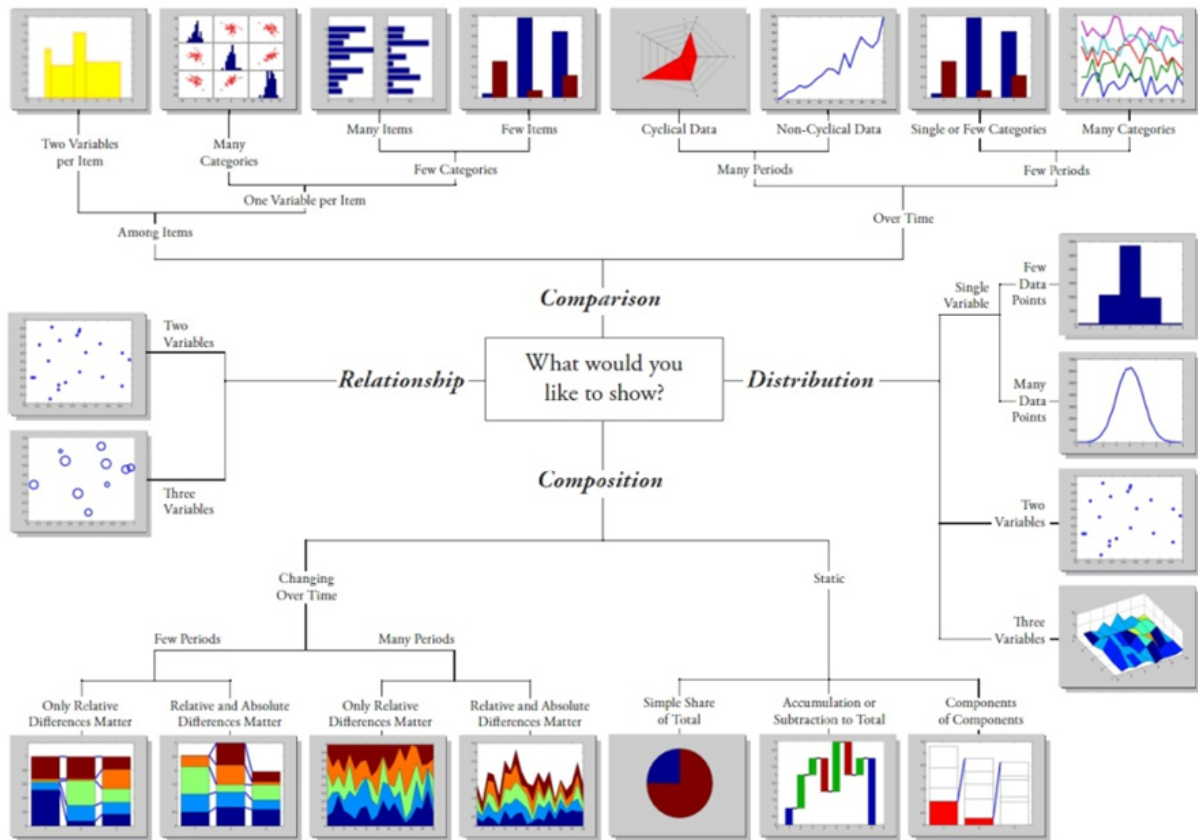


Figure 9.2: An alternative way to group chart types [D. Hull et al., A. Abela et al.].

The take home point here is that while it's possible to group data visualization types in a number of ways, it is important to hone our own sense of the most appropriate visualizations for particular situations, which may be informed by schema developed previously, and so on (see Figures 9.1 and 9.2).

As we have already discussed in Chapter 2, regardless of our preferred **organizing principles**, there are some data visualizations that are true **workhorses** – they are usable in some way with almost any dataset, will be familiar to most lay people and are particularly useful for exploring or presenting data. Other visualization types have more **situation-specific uses**, and may be difficult to use or ineffective in other situations – practice and discernment are required to use them **skillfully** and **appropriately**.

Different people may consider different options for a list of workhorse visualizations. For instance, we could trot out the following “old-faithfuls”:²

Data Exploration

- Scatterplots
- Rug charts
- Histograms
- Two-way tables

Data Presentation

- Text blocks
- Tables
- Line graphs
- Bar charts

These categorizations – **exploration** vs. **presentation** – are not intended to be hard and fast: at times, we might use a barchart for exploration and a scatterplot to present a key finding, say. Everyone will develop their own approach for the use of these visualizations, but we provide some comments and pointers as guidelines.³

9.1 Exploration Visualizations

Scatterplots

A **scatterplot** is one of the fundamental tools of data scientists. In a rapid and accessible manner, it can reveal **key relationships** between two variables simply by plotting the data points on a **grid**.⁴

When using a scatterplot for **data exploration**, the emphasis lies not only on the aesthetics of the plot, but also on what a simple version of a plot show about the **nature of the relationship** of the variables being plotted – in particular, the **emergent patterns**. Do a circular cloud of points, a straight line, a diagonal line, a wavy line appear? All of these different patterns denote different potential relationships between two variables which could be worth exploring further.

2: Another selection has already been reviewed in Chapter 2.

3: As a rule of thumb, these lists provide a good starting point when deciding which tools to pull from the toolbox first.

4: Which we technically refer to as a Cartesian plane, a two dimensional cartesian coordinate system.

But caution must be exercised when using scatterplots for **data presentation**. Because they can represent all of the points in a dataset, there is a risk for **clutter** (and overwhelming the consumer). The message can easily get lost. Consequently, we suggest only using scatterplots for communication when the pattern is naturally clear and relevant to the broader context of the story being presented.

Bubble charts, which are a variation on scatterplots, can communicate relationships between **multiple dimensions** and provide a powerful tool to render multivariate relationships. We will discuss them further in a coming section of this chapter.

Scatterplots (and bubblecharts) are most commonly used for **quantitative data**, but it is possible to have one or both of the axes represent a qualitative variable, using an approach similar to that taken on the horizontal axis (x-axis) of bar charts. This can lend itself to misinterpretation, however.

Lastly, it is not uncommon for scatterplots to be overlaid with a **trend line or curve**, such as in the data storytelling tropes of Section 8.3 (*Evolving a Storytelling Chart*).⁵

5: Creating trend lines involves calculations over the data points represented; this can be automated by the software used to render the chart, as we will see in the next few chapters.

TL;DR Summary and Comments:

- plots show relationship between 2 variables (scatterplot) or 3 variables (bubble plot)
- we can use average lines (or similar curves) to provide context
- consider using groupings to add clarity (e.g., colour gradients)
- colour and geometry allow us to plot (at least) 2 extra variables on a 2D scatterplot
- the data may need to be re-scaled or binned
- a movie could be used to visualize an additional ordinal variable
- text can also be added to visualize an additional categorical variable
- works best when chart is not too encumbered

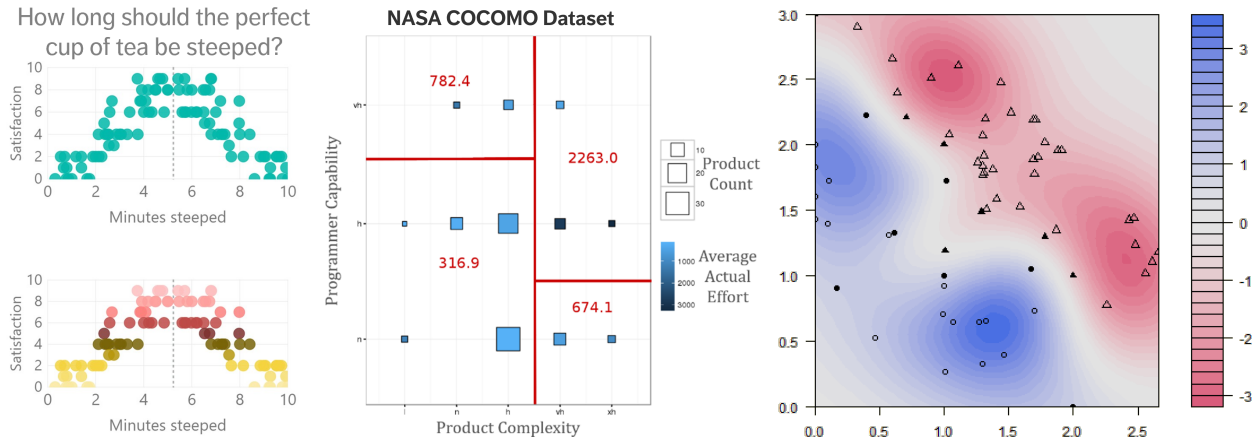
Examples of scatterplots (and bubble charts) are provided in Figures 9.3 and 2.7.

Rug Charts and Histograms

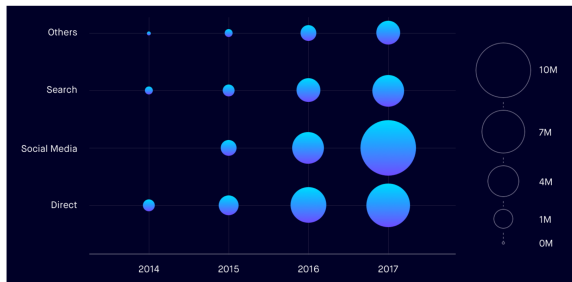
A **rug chart** is essentially a single quantitative variable that is plotted on a horizontal or vertical **number line** (see Figure 2.4). It is easily overlooked as a strategy to understand the behaviour of each variable separately from any of the other variables in the dataset. In some ways it can be thought of as a “quick and dirty histogram”, since histograms also focus on a single variable, albeit in a slightly more nuanced fashion.⁶

6: Rug charts are useful building blocks for more complex visualizations like **radar** or **spider charts**.

Among data visualization approaches, **histograms** are at once one of the most familiar and one of the least well-understood. It is invaluable for data exploration, but can be treacherous when used in data presentations, to the point that to the question of when a histogram should be used in a data presentation context, our answer is: **never, probably**.



Website Traffic from Different Channels Over a Four-Year Period



Wine Quality Relative to Three Factors

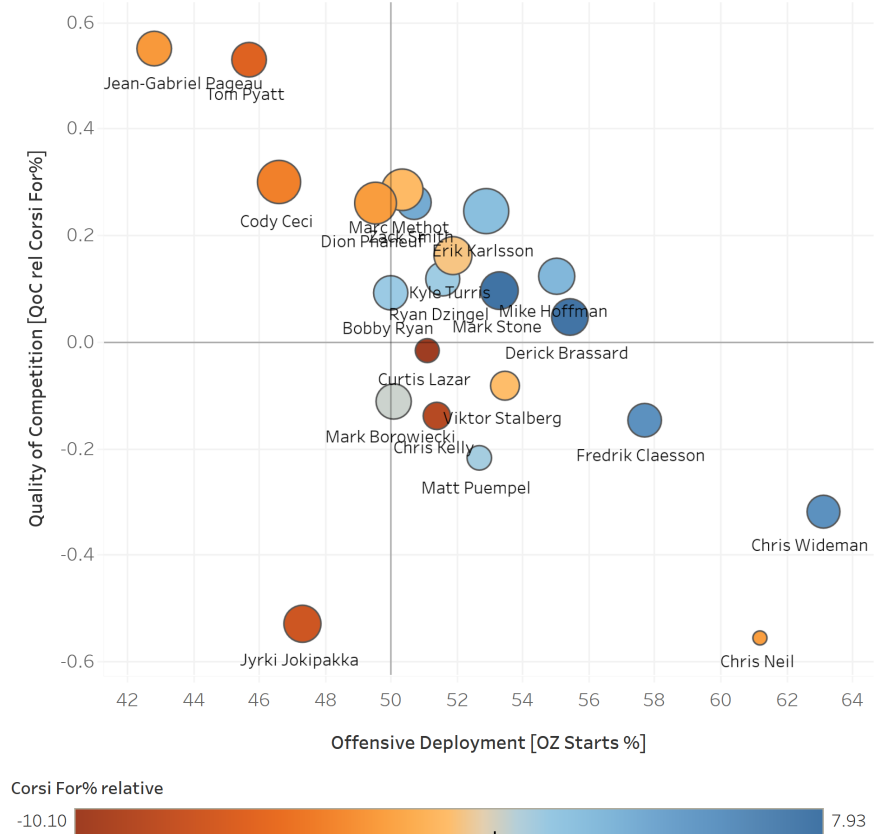
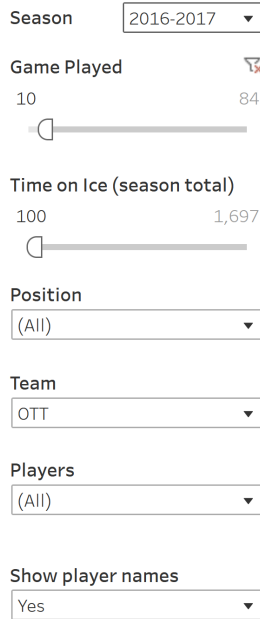
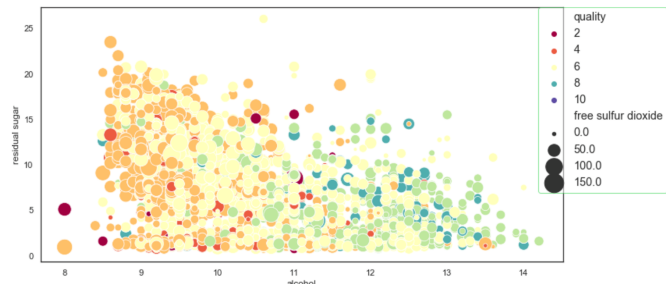


Figure 9.3: Scatterplots and bubble charts: personal collection (top row); Medium [↗](#) (middle left); Towards Data Science [↗](#) (middle right); Ottawa Senators player usage, 2016-2017 ([Hockey Abstract ↗](#), bottom row).

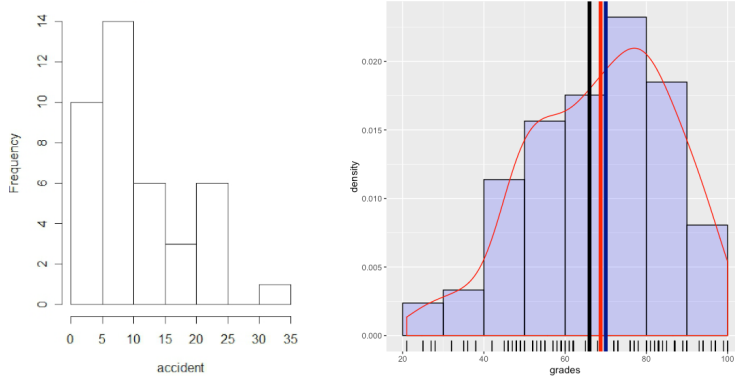
As with the number line, the histogram focuses on a single **quantitative variable**. Once it has been selected, we must then carry out some calculations on this variable. If $\{x_i \mid i = 1, \dots, n\}$ are the values taken by the variable in the dataset, then its histogram should contain the following information:

- the range of the histogram is $r = \max\{x_i\} - \min\{x_i\}$;
- the number of bins should be $\approx k = \sqrt{n}$, where n is the sample size;
- the bin width should approach r/k ,
- and the frequency of observations in each bin is then used to determine the height of the bar.

By abstracting away from the “bumpiness” of the individual bars, it is possible to get a general sense of the abstract variable shape in the dataset.⁷ Examples of histograms and rug charts are shown in Figures 9.4, 2.5, and 2.6.

7: Histograms are related to **bar charts**, which we will discuss shortly.

Figure 9.4: Examples of histograms and rug charts: frequency of daily number of road accidents in Sydney, Australia, over a 40-day period (left); grade distribution in a 2nd year probability and statistics class, with mode (black), median (blue), and mean (blue) overlay, density curve, and rug chart (right).



Two-Way Tables

Most of our discussion on exploration visualizations has focused on quantitative variables; a straightforward preliminary strategy for exploring **qualitative variables** is through the use of a **two-way table**,⁸ which displays the counts of one variable level relative to the counts of a second variable level, as is illustrated in Figure 9.5.

8: Or m -way table, more generally.

Figure 9.5: Example of a two-way table for an artificial dataset with 89 observations and 2 variables: window type and window size; for example, there are 11 medium-sized door windows among the 80 observations.

	Large	Medium	Small
Window	1	32	31
Door	14	11	0

We could use $\binom{n}{2}$ pair-wise 2-way tables to display information for a dataset with $m = 3$ categorical variables (speed, size, season), there would thus be three 2-way tables: speed \times size, speed \times season, and season \times size, say.

	large	medium	small		autumn	spring	summer	winter		large	medium	small
high	13	56	73	high	32	34	38	38	autumn	19	33	28
low	32	24	2	low	16	13	12	17	spring	21	34	29
medium	38	56	46	medium	32	37	36	35	summer	19	36	31
									winter	24	33	33

But this only presents a part of the picture. The 1-way tables provide a **univariate summary**:

season	autumn	spring	summer	winter
	80	84	86	90
size	large	medium	small	
	83	136	121	
speed	medium	high	low	
	140	142	58	

while the 3-way table speed × season × size provides the full picture, but it is not the only way to represent it:⁹

9: What would the other combinations (size × speed × season, etc.) look like?

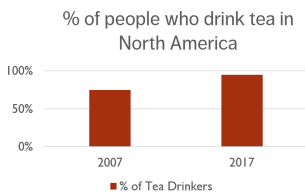
	speed = high			speed = low			speed = medium		
	large	medium	small	large	medium	small	large	medium	small
autumn	3	13	16	9	6	1	7	14	11
spring	3	14	17	7	6	0	11	14	12
summer	3	16	19	6	6	0	10	14	12
winter	4	13	21	10	6	1	10	14	11

9.2 Presentation Visualizations

Text Blocks

The simplest of the presentation visualizations, the **text block**, may not even seem like it belongs in a list of data visualizations, given its lack of visual elements outside of the written word. However, when treated as a graphical element, where the focus is on a fact containing one or two numbers at most, they are excellent at “**setting the scene**”.

This is particularly useful in a dashboard or in a report context, where text blocks can be used to draw the focus to an area of the report which contains a more detailed breakdown or analysis of the data in question.



95% of the population
drinks tea today compared to
75% in 2007

Figure 9.6: Bar chart (left) vs. text block (right): in this case, the chart is overkill – the insight is much more easily conveyed with text.

Don't neglect this simple approach to conveying insights.

Tables

Tables are another text-heavy visualization which interact with our **verbal system**: we **read them**. They are useful for comparing values across variables.

One complicated aspect of tables is that the audience has considerably leeway in regard to **how they elect to read** them: they may focus on the relationship between numbers **across each row**, or **down each column**. Furthermore, if the data relates to specific individuals (or units) audiences are expected to be most interested in, they will certainly look for and focus on their rows, potentially to the exclusion of everything else.¹⁰

Importantly, table design should **blend into the background**.¹¹ It is the data that should stand out, not the borders – if you must display large, dense tables,¹² consider alternating the table row colour from white to a very lightly noticeable shade to help the eye scan across the rows and separate one row from another.

Name	Last Year	This Year	Name	Last Year	This Year
Bob	20	30	Bob	20	30
Fred	30	40	Fred	30	40
George	10	15	George	10	15

- 10: Designers may need to use the Gestalt guidelines of Section 4.2 to draw the audience’s eye to another location in the table.
- 11: Although it seems to go against easily-accessible MS PowerPoint templates...
- 12: Must you, really?

Figure 9.7: Fanciful table (left, not recommended) vs. simple table (right, recommended).

The **table heat map** provides a variant on the table, where cells contain a colour as well as a number, and in which the colour is leveraged to convey **magnitude**, by mapping the colour hue and saturation to the cell value.¹³ Eventually, the numerical values or cell text labels may be removed **without altering the message**, leading to a more holistic reading of the data visualization.

- 13: A **single colour saturation** with a legend (white = low, blue = high) is preferable to **colour differentiation** (rainbow scale).

	Last Year	This Year	Next Year	Optimum		Last Year	This Year	Next Year	Optimum		Last Year	This Year	Next Year	Optimum
George	20	20	20	20	George	20	20	20	20	George				
Peter	40	35	30	25	Peter	40	35	30	25	Peter				
John	10	10	5	5	John	10	10	5	5	John				
Sandra	25	30	35	40	Sandra	25	30	35	40	Sandra				

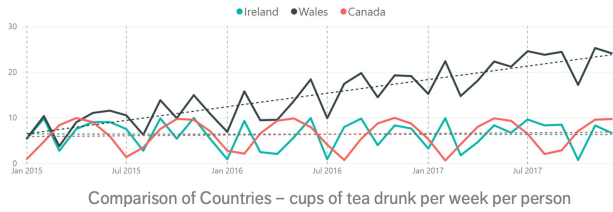
Figure 9.8: From table (left) to heat map table (middle) to holistic heat map table (right).

Line Graphs

Line graphs have some similarities with scatterplots. As a consequence, it can be difficult for people new to data visualization to appreciate how they differ from the latter and to get a good sense of when to choose one over the other. Both involve plotting data points on a grid, and both can accommodate curve overlays, for instance – the critical distinction is that in the case of a line graph, the line (curve) connects **all the points in sequence**,¹⁴ and passes **exactly once through each** of the points.

While some line graphs only display one line (curve), they can also be used when the dataset contains also at least one categorical variable whose levels

- 14: The dataset must thus contain an ordinal variable that is used to create the sequence of points, and typically at least one quantitative value used to place the point on the chart.



	Start	Monthly Number of Cases	End	Low	High	Mean	Std Dev	Blanks	Zeros	Trend
TOTAL	19502		17265	15150	25072	19903	2612	0.0	0.0	379.2
Hospital #1	46		19	3	46	19	9	0.0	0.0	-1.6
Hospital #2	156		240	101	326	194	60	0.0	0.0	9.7
Hospital #3	16		11	2	76	15	15	0.0	0.0	-2.9
Hospital #4	3		13	0	105	9	15	0.0	0.4	-1.8
Hospital #5	42		50	25	91	61	16	0.0	0.0	1.2
Hospital #6	48		53	34	169	67	25	0.0	0.0	0.6
Hospital #7	0		N.A.	0	0	0	0	2.2	9.8	0.0

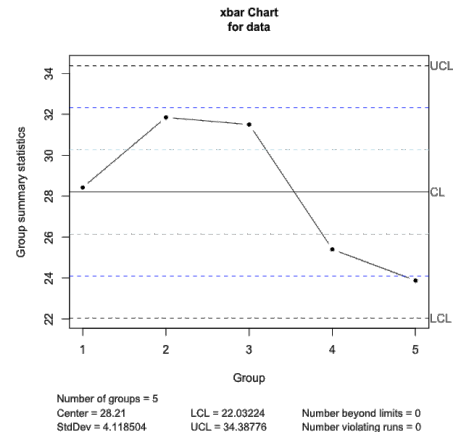
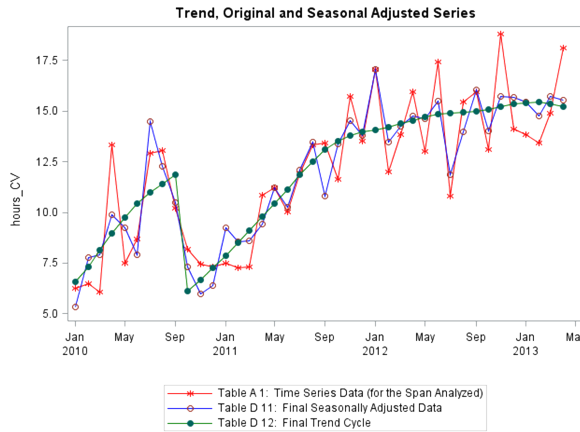


Figure 9.9: Various line graphs and sparklines, from personal files.

can be used to separate the observations in different groups. Each of these categories can then be given its own curve on the line graph, allowing us to compare the behaviour of the qualitative values across the categories.

Although describing the mechanics of a line chart can be complicated, they are typically familiar to a broad audience and so are likely to be more readily **interpretable** than some other data visualizations (as in Figure 9.9).¹⁵

TL;DR Summary and Comments:

- line charts can show a single series or multiple series of data;
- they particularly useful to display time series;
- axis scale should be clear and relevant;
- the *y*-axis should be “anchored” when using dynamic filters so that the graph does not “jump around” as users interact with it.

Bar Charts

One of us¹⁶ thinks that the **bar chart** is the workhorse of all workhorse among data visualizations. It is almost immediately familiar to most people and, if used in an expected fashion, **readily interpretable** as well. Although some people may hesitate to use a bar chart simply because it is so familiar and frequently used, we believe that this is a strength, a weakness.¹⁷

The **basic bar chart** represents a single numeric variable broken down by the values of a categorical variable.¹⁸

These charts are quite versatile and useful. Apart from very rare instances, they should always have a **zero baseline**. When constructing them, we

15: **Sparklines** are a line graph variant which are meant to play the role of a “word” in a larger chart or paragraph [3].

16: *cough* Jen *cough*

17: If novelty is desired, there are many variations on the bar chart, both with respect to what types of data are incorporated and aesthetic and presentation choices, that can add interest and nuance to the basic bar chart.

18: See Section 2.2 for more information.

19: Horizontal charts are apparently easier to read, as we have discussed in Section 4 [38].

20: Variations include: pie charts, gauge charts, funnel charts, lollipop charts, waterfalls, stacked bar charts, cluster bar charts, 100% bar charts, percentage bar charts, etc.

21: Although that is not always the case.

22: Not that there is anything wrong with that, of course.

recommend using either the **graph axis** or the **data labels**: axis for broad statements, data labels for detail information.¹⁹

From a design point of view, the basic bar chart can be transformed in a number of ways which impact how viewers interpret and prioritize different aspects of the visualization.²⁰

Funnel charts are typically used to represent decreasing proportions amounting to a 100% total.²¹ These can be very useful to help audience **quickly prioritize** items without having to actively filter the data.

Gauge charts are often used as a dashboard component (with or without needle) – they typically display single value measures on the way to some goal or **key performance indicator (KPI)**, in a manner that can quickly be scanned and understood. While gauge charts are particularly useful to show progress, they may ultimately prove to be a management fad.²²

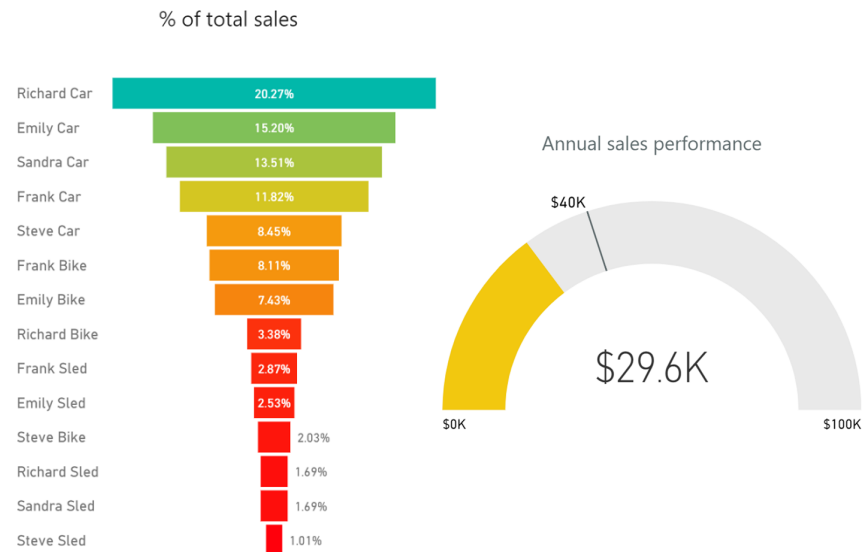


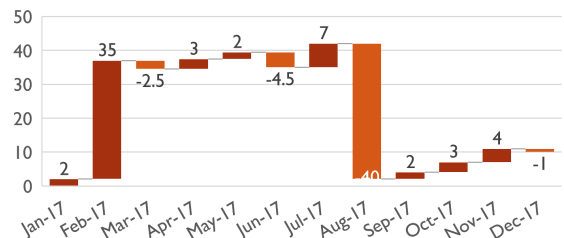
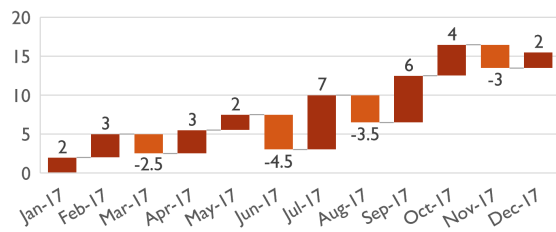
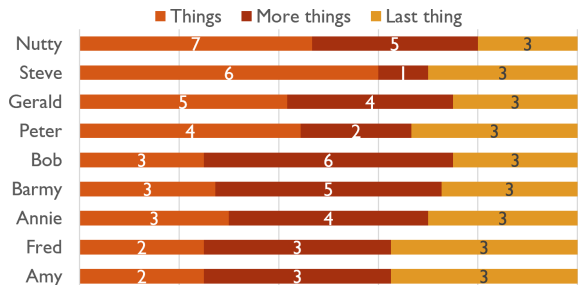
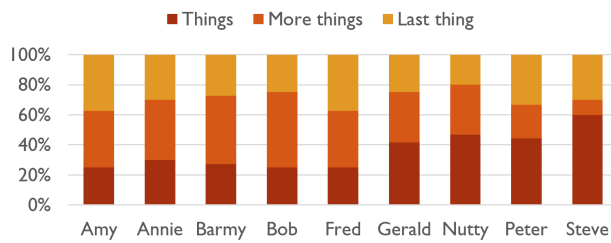
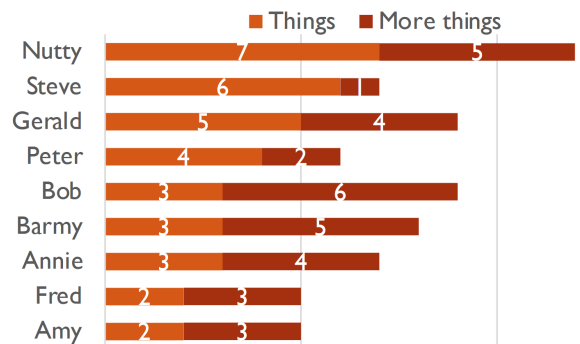
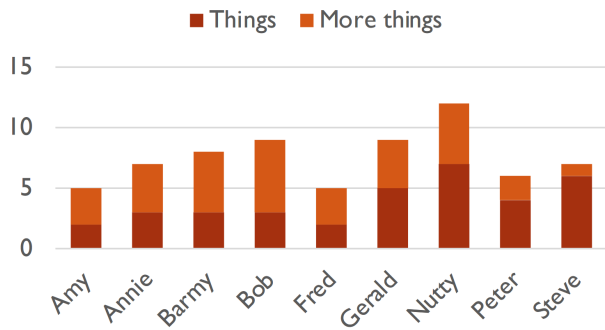
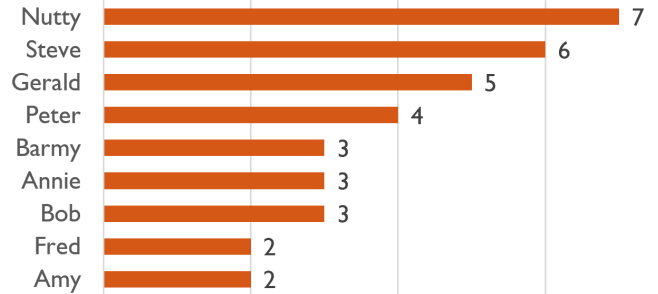
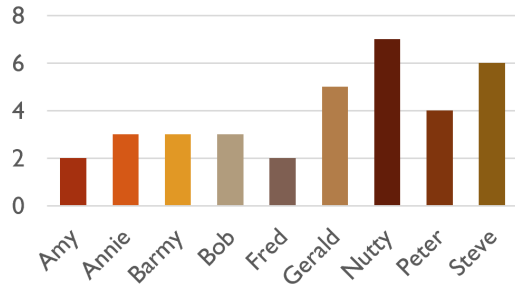
Figure 9.10: A funnel chart representing the % of total sales for each salesperson × product (fictional example, left); gauge chart, with target (right).

Stacked bar charts are designed for comparing totals, but can quickly become overwhelming. They are hard to sort and order. Filtering is complicated in dashboard applications like Power BI because it is unclear how the chart should respond when filter is applied.

100% bar charts work well for visualizing portions of a whole on scale from negative to positive. They have a consistent baseline at each of the extremities (either left/right, or top/bottom), making it easy to compare the bars. The issue, however is that there is no relative measure of the magnitude of data. As with other bar charts, research shows that horizontal is easier to process than vertical.

Waterfall charts shows how the initial value increases or decreases using a series of intermediate values; different colours should be used to represent increases and decreases. One drawback is that it is difficult to remove charts elements without removing context.²³ Note that large increases or decreases may look odd (as in Figure 9.11).

23: In other words, it is difficult to declutter waterfall charts.



Number of Units Sold

Number of Units Sold

Figure 9.11: Bar charts and variants: basic (top row); stacked bar charts (2nd row); 100% bar charts (3rd row); waterfall charts (bottom row).

9.3 The Rest of the Visualization Landscape

As we have seen throughout (especially in Chapter 3), modern data visualization endeavours often go above-and-beyond the workhorse visualizations. In this section, we will present some of the more sophisticated approaches used in data presentations.

Maps, Heat Maps and Choropleths

A more thorough treatment of **geographical maps** and **map-based** visualizations would probably require a chapter in its own right. Most of us are quite familiar with geographical maps, so they tend to be easier to interpret; we can play with this to produce a striking effect when the data visualization shows unexpected results, and thus change the viewer’s perception in the process (see Figure 9.12).

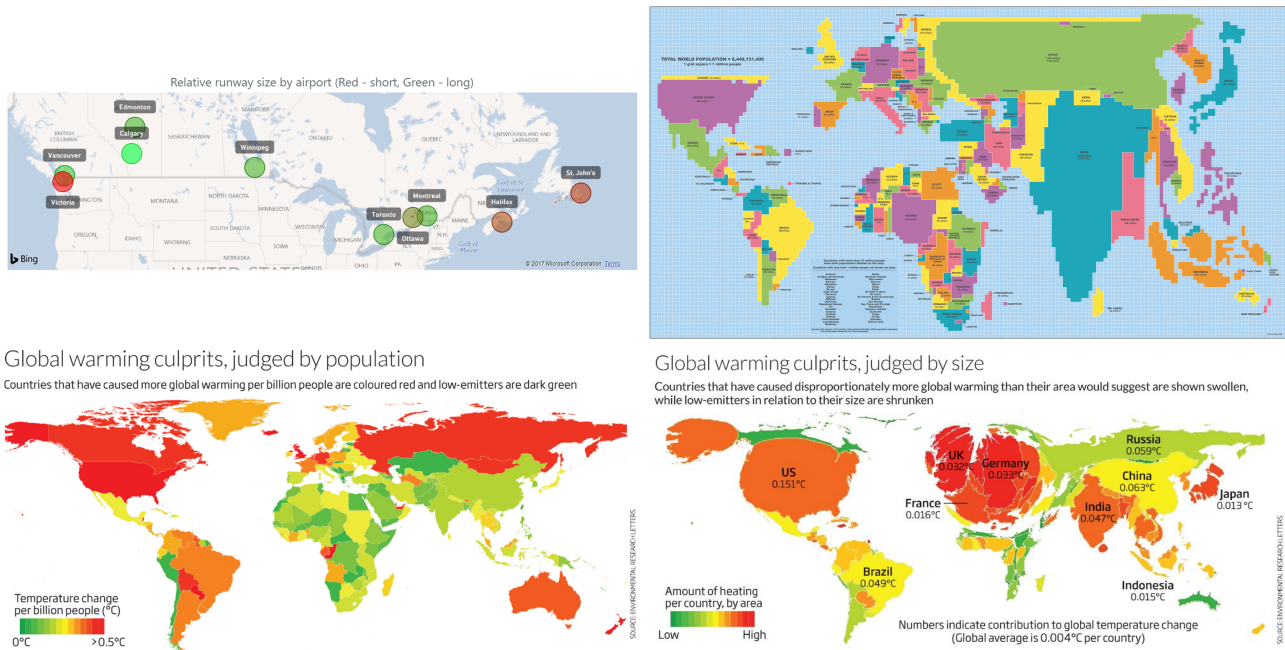


Figure 9.12: Maps, maps, maps: a sprinkle of maps and distortions – Canadian airports (top left, personal file); population cartogram (top right, Paul Breding); global warming culprits, by population and by size (bottom row, New Scientist).

Heat maps are ideal when we want to look at the relationship between 3 or 4 variables. If one of these represents a percentage or a value within a set range, it can be used to fix the colour scale, for comparison purposes. The other variables are then used to locate and size markers on the display.

If the axes variables are continuous, it could still be preferably to **bin them**: this decreases the number of required observations for usefulness. It is typically easier to read such charts if colours are selected along **natural colour gradients**, such as White → Blue or Red → Black.²⁴ When the background canvas is a non-distorted geographical map, heat maps are known as **choropleths** (see Figure 9.13).

24: More sophisticated gradients can be used (Red → Yellow → Green, say), but those are less than ideal from a Gestalt perspective, or if some viewers are colour blind – this is another clear case where “less is more”.

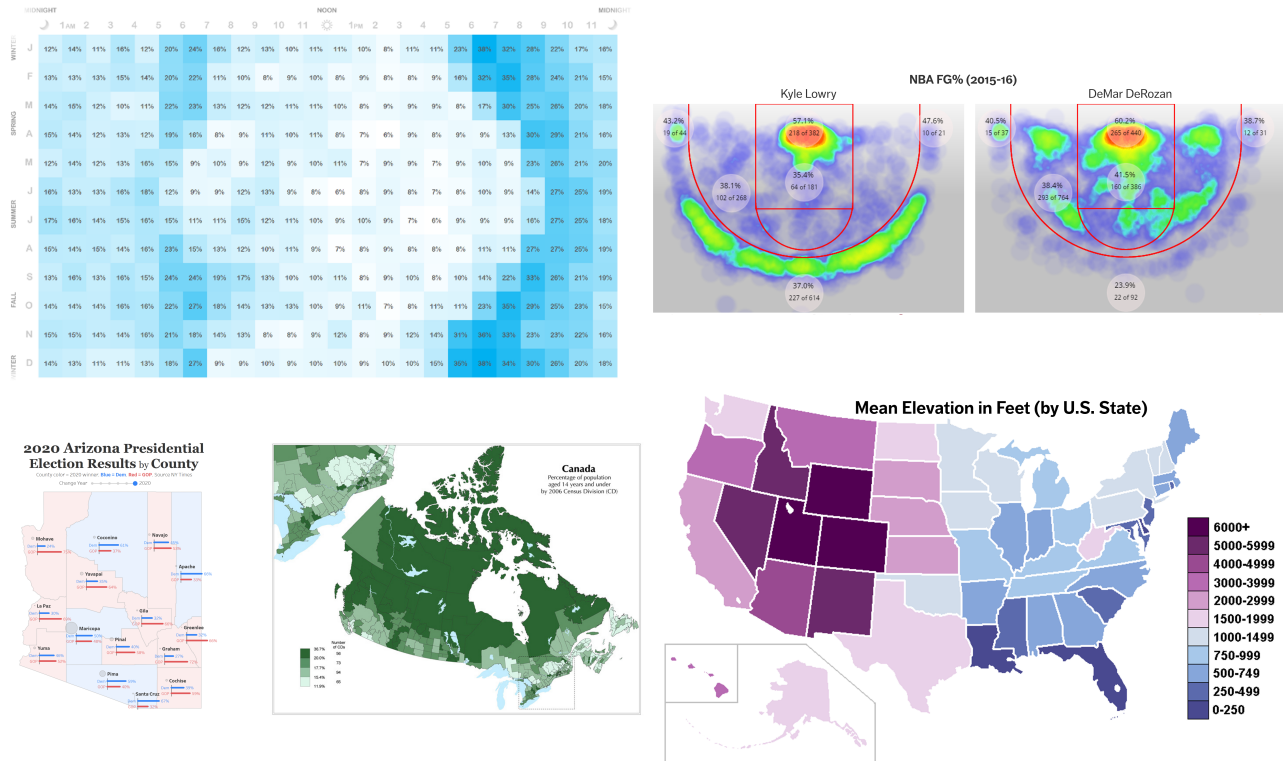


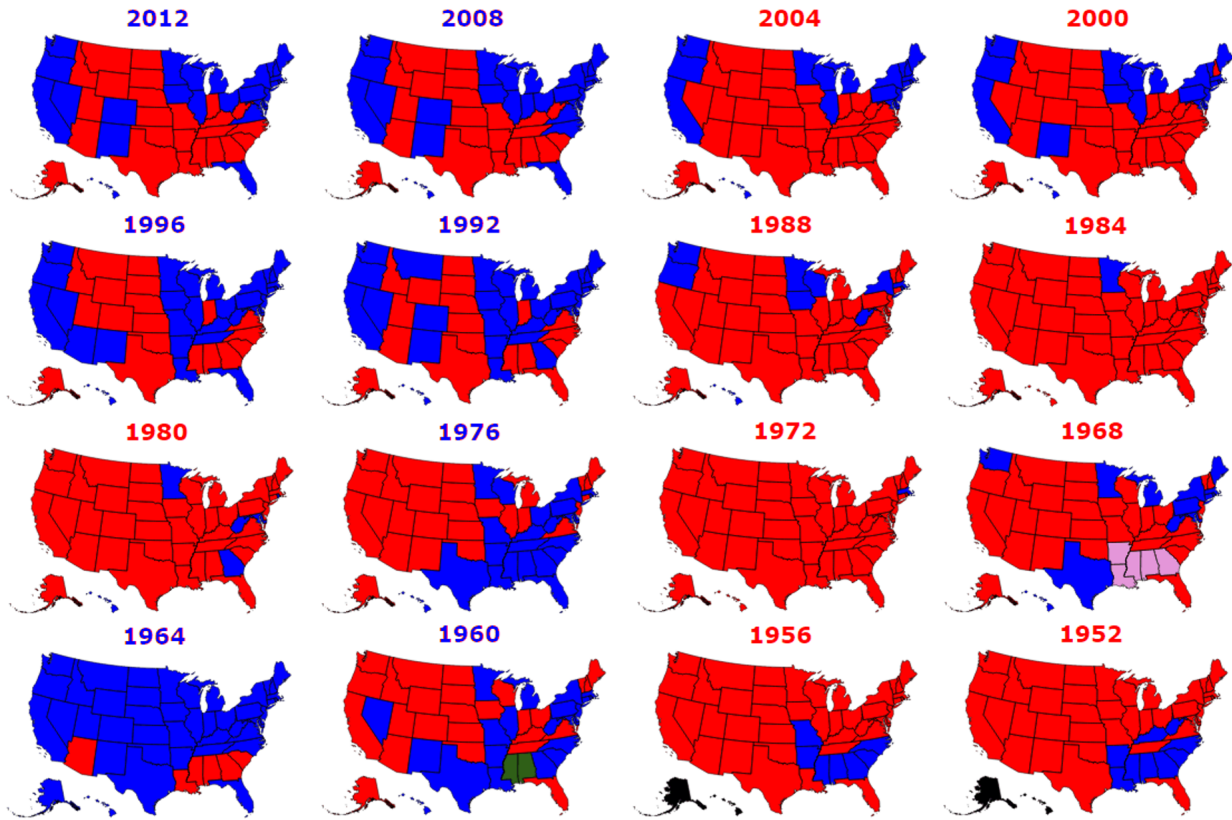
Figure 9.13: Heat maps and choropleths: The Horizon or Pedestrian Risk (J. Nelson, IDV Solutions, top left); basketball shooting charts (NBAsavant.com, top right); Election choropleth (A.E. McCann, bottom left); Canadian population choropleth (Statistics Canada, bottom middle); US elevation choropleth (author unknown, bottom right).

Bubble Charts

Unlike scatterplots, which have already been discussed both in Chapter 2 and earlier in this chapter, **bubble charts** can serve to illustrate the interactions between multiple variables (when used correctly). Importantly, however, they are usually most useful when there relationships between the variables in questions are **strong**, resulting in **clear patterns** in the chart. We must also be careful when choosing how to represent the many variables involved – there are likely more bad options than good choices on this front, and so experimentation is likely to be required; examples can be seen in Figures 1.4 and 9.3.

Small Multiples

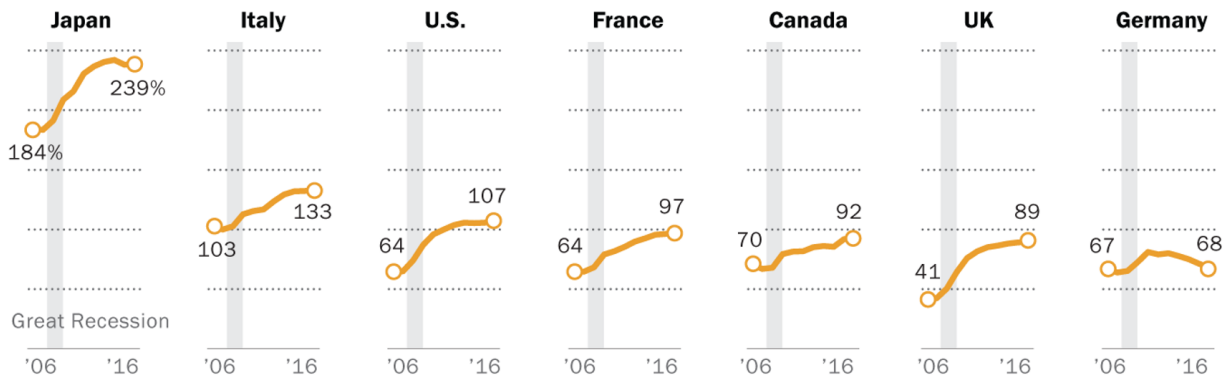
Combining multiple visualizations of the same type and presenting them as tiles in a larger composite visualization can have a powerful effect – we call such combined charts **small multiples**, after Tufte [3]. Depending on the choice of visualization, the small multiple can depict change over time, as in a flip book, or encourage high-level comparisons across categories (see Figure 9.14).



U.S. Electoral College Results 1952 – 2012

After Great Recession, debt increased substantially in most G-7 economies

Total gross debt as a share of GDP in the Group of Seven nations



Note: Gross debt represents total liabilities of all levels and units of government – national, state/provincial and local – less liabilities held by other levels or units of government, unless otherwise noted by source.

Source: The International Monetary Fund, World Economic Outlook, accessed Sept 7, 2017.

PEW RESEARCH CENTER

Figure 9.14: Small multiples: US electoral results choropleths, by year (author unknown, top); debt line graphs, by G7 country (Pew Research Center, bottom).

Area Charts and Treemaps

Area charts use physical areas to represent various quantities, such as in Figure 9.15.

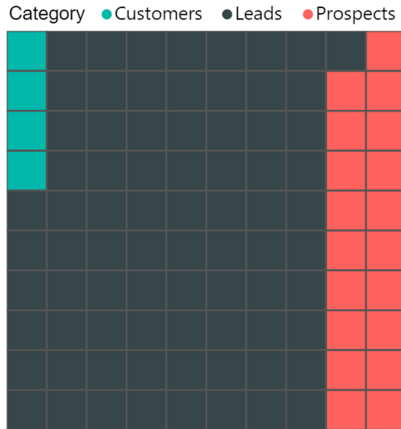


Figure 9.15: Area chart with three categories.

We suggest, however, to try to avoid area charts, except in situations where the plotted quantities have **vastly different magnitudes** as human brains have a hard time attributing a value to a 2D area (see Section 2.4).²⁵

An exception to this warning might be the **treemap**, if used with care. A treemap can simultaneously show the big picture and compare categories (or sub-categories) easily. They are useful for prioritizing “big ticket items” in dynamic dashboards.²⁶

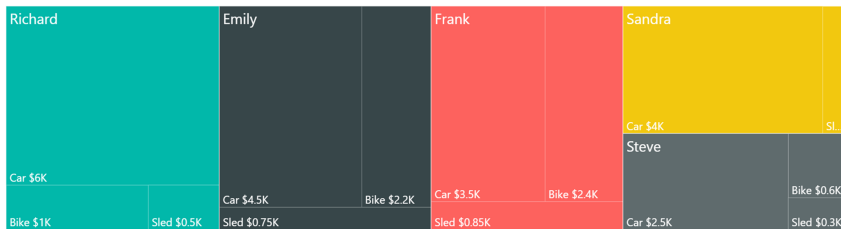


Figure 9.16: Treemap with five categories and three sub-categories.

Text Visualizations

Not to be confused with text blocks, **text visualizations** use text attributes (such as size and colour) to represent some other variable associated with the words. For maximal impact, font size may be a function of frequency. These visualizations are typically used for univariate categorical data, but small multiples, cloud shape, word placement, colour, and hue could be used to integrate more variates.

In many implementations, the word placement and colour choice algorithms are “hidden” from the users. As an example use case, text visualizations can be used to answer **authorship questions**.

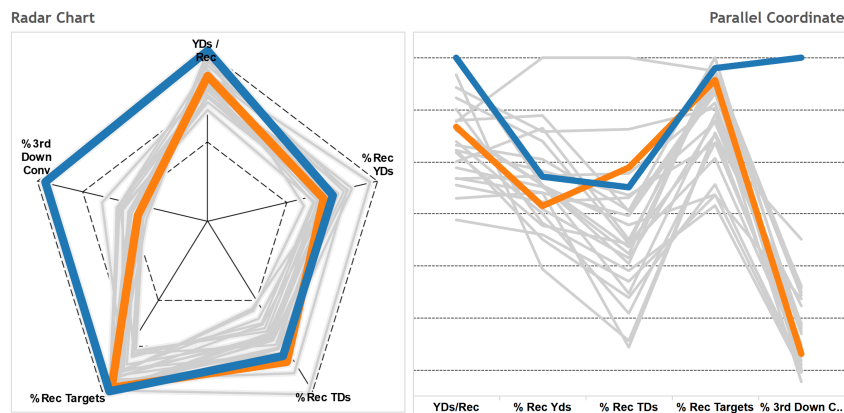
25: Area maps (such as a pie chart) are perhaps most useful when viewer need to see that a quantity is much greater than another, but the numerical factor by which it is not relevant.

26: Although labeling and colouring can be tricky... but that is potentially a problem with just about any data visualization.

Radar Charts and Parallel Coordinates

Although **parallel coordinate charts**, which stack and connect multiple **rug charts** to show relationships between potentially large numbers of variables, are a relatively obscure type of visualization, a variation has been increasing in popularity in recent years.

Radar charts, which arrange the axes radially as spokes coming out of a central point, are often seen in social science or business contexts, where they are used to show survey results. When used in this manner, the overall shape of the connected line on the radar chart gives a gestalt sense of the response profiles.²⁷



27: E.g., are they mostly low or mostly high? and so on.

Figure 9.18: Radar chart and parallel coordinate chart of NFL players Cruz vs. Fitzgerald performance (A.E. McCann).

Trees and Networks

Using **networks** to both model and visualize systems can give us insights into the system. Having a solid conceptual understanding of the system through the use of these visualization types can help us draw legitimate and sound conclusions.²⁸ Examples are provided in Figure 9.19.

Animated and Interactive Visualizations

Animation and **interactivity** do not always improve a visualization. What insights can they provide? That depends on the data, and on the visualization method.

Even when done well, 85% of users don't bother with interactive viz, according to a [NY Times](#) analysis of their own visualizations at *The Upshot*. This very strongly supports the notion that the **default visualization** (i.e., the one that greets viewers when they first load the website on which it is found) should be coherent and self-consistent as is (see Figures 9.20 and 9.21 for examples).

28: Here are some clues that suggest that using a tree or network visualization could be useful:

- are we dealing with flow (of something) along pathways?
- are we dealing with a collection of objects that input and output things?
- are the inputting/outputting objects homogeneous?
- are we dealing with relationships and connections between objects?
- Are we dealing with a situation where one object influences another object?

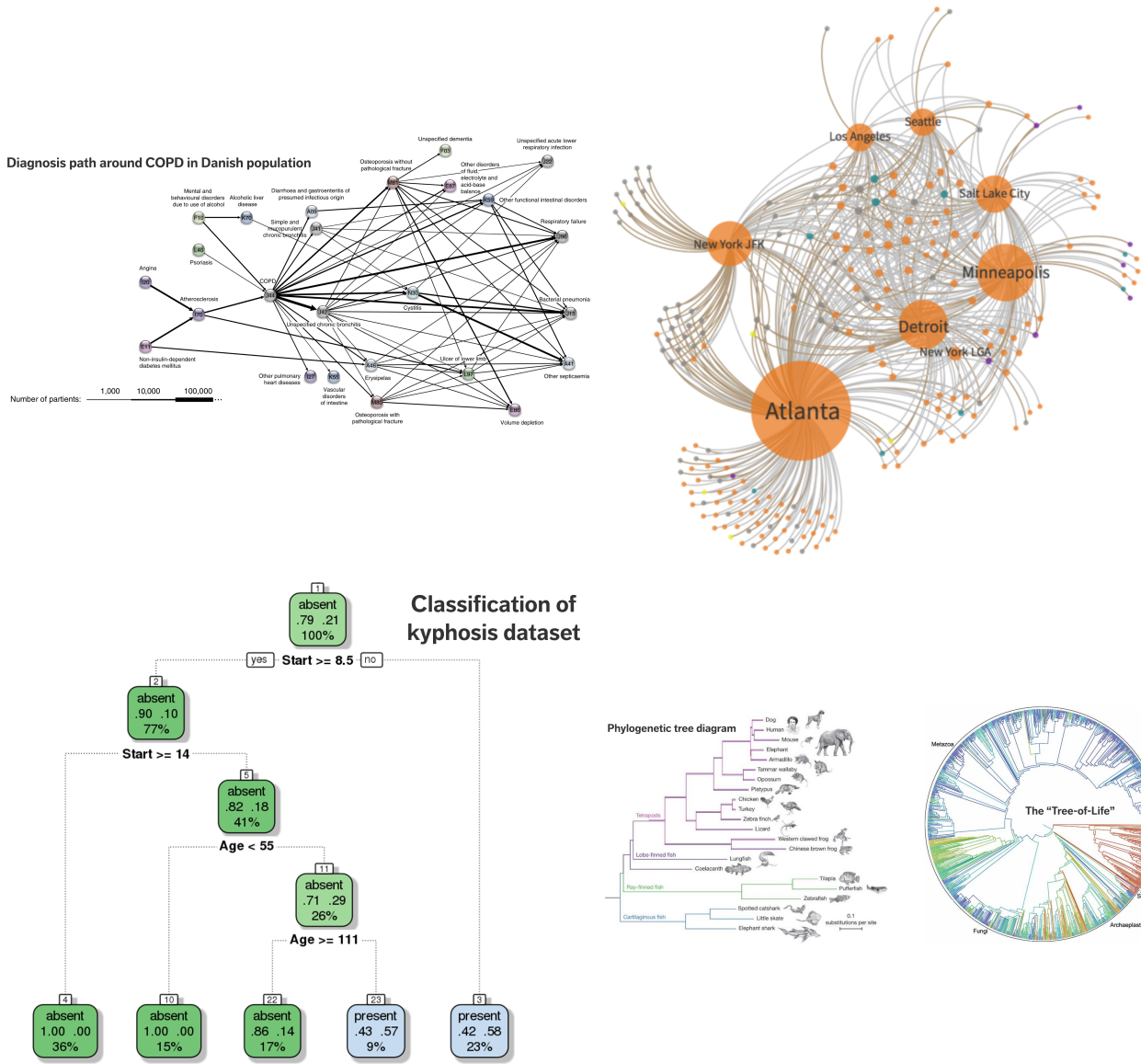


Figure 9.19: Trees and network diagrams: disease progression [32] (top left); US airport hubs (top right, author unknown); classification [1] (bottom left); tree of life [2] (P.Z. Meyers, bottom right).

Brazil 2014

THE CLUBS THAT CONNECT THE WORLD CUP

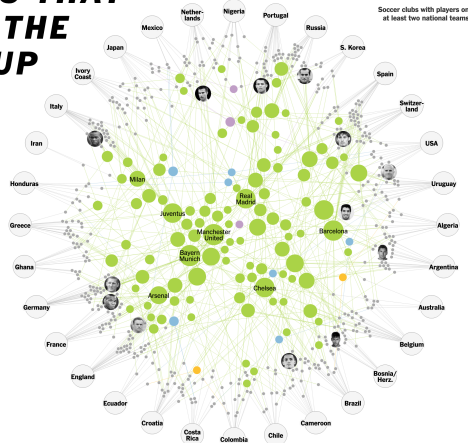
BY ORCOUR ANDER JUN 20, 2014

The best national teams come together every four years, but the global tournament is mostly a remix of the professional leagues that are in season most of the time. Three out of every four World Cup players play in Europe, and the top clubs like Barcelona, Bayern Munich and Manchester United have players from one end of the globe to the other.

- Europe
- Africa
- Asia
- South America
- North America

Brazil vs. Argentina

Even archrivals Brazil and Argentina overlap. Neymar, Brazil's star forward, plays alongside Lionel Messi, the Argentine captain, on powerhouse Barcelona. In all, eight Brazilians and 22 Argentines play together on European club teams.



This Chart Shows Who Marries CEOs, Doctors, Chefs and Janitors

By Adam Pearce and Dorothy Gambrell February 11, 2016

When it comes to falling in love, it's not just fate that brings people together—sometimes it's their jobs. We scanned data from the U.S. Census Bureau's 2014 American Community Survey—which covers 85 million households—to find out how people are pairing up. Some of the matches seemed practical (the most common marriage is between **grade-school teachers**), and others had us questioning Cupid's aim (why do female **dancers** have a thing for male **engineers**). High-earning women (**doctors**, **lawyers**) tend to pair up with their economic equals, while middle- and lower-tier women often marry up. In other words, female **CEOs** tend to marry other CEOs; male CEOs are OK marrying their secretaries.

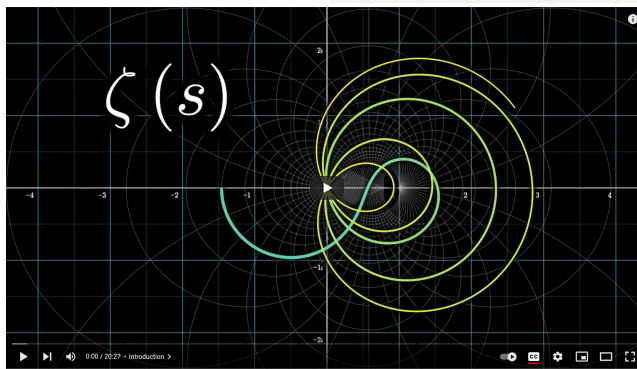
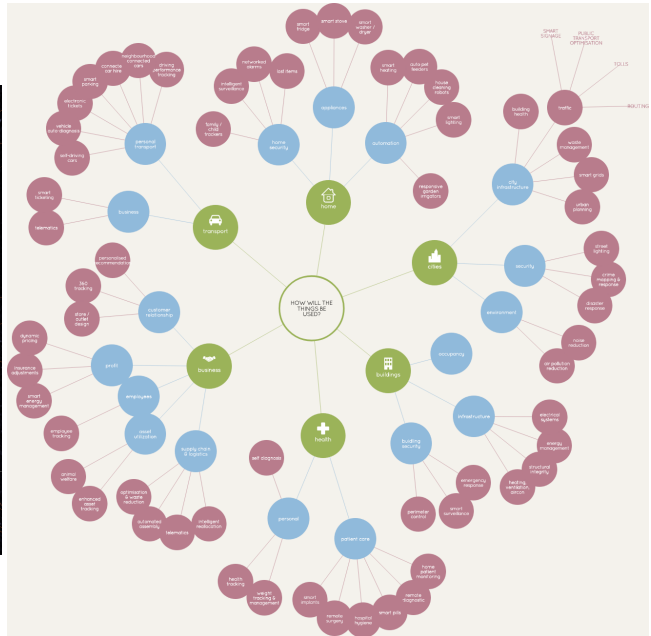
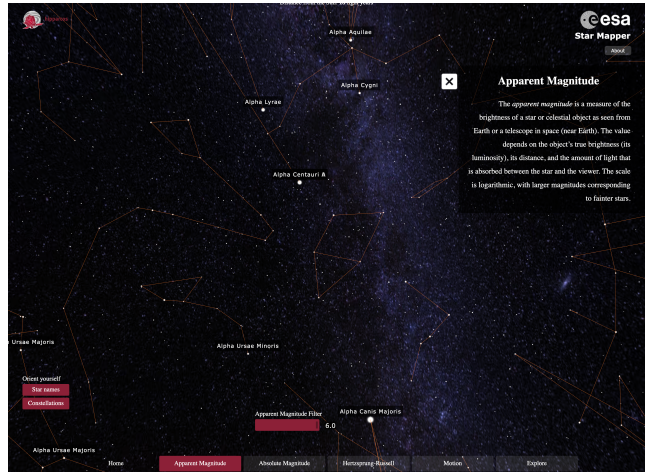
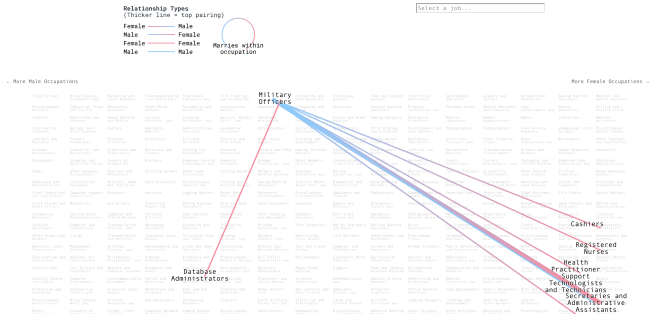


Figure 9.20: Animated and interactive charts: *The Clubs That Connect the World Cup* , NY Times, 2014 (top left); *Who Marries Whom* , Bloomberg, 2016 (top right); *Hipparcos Star Mapper* , European Space Agency, 2016 (middle left); *The Internet of Things – a Primer* , Information is Beautiful, 2016 (middle right); *Visualizing the Riemann ζ Function and Analytic Continuation* , 3Blue1Brown, 2016 (bottom row).

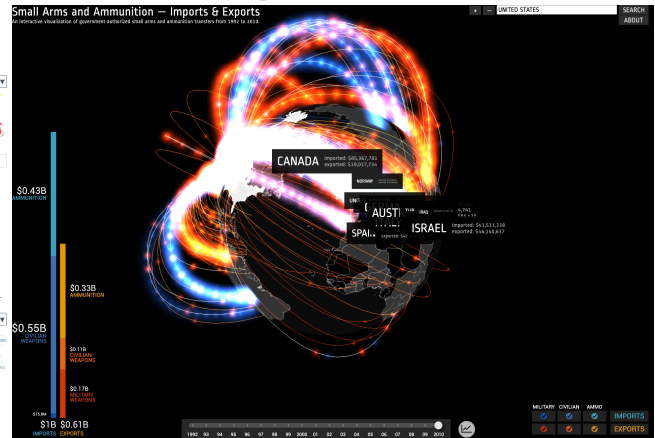
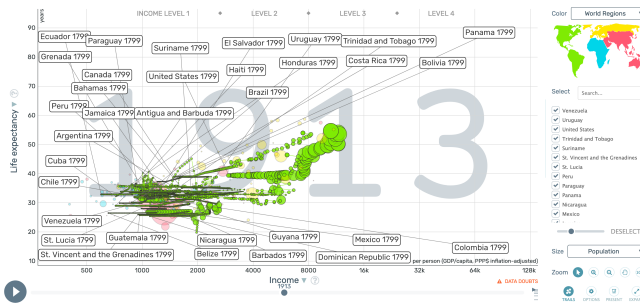
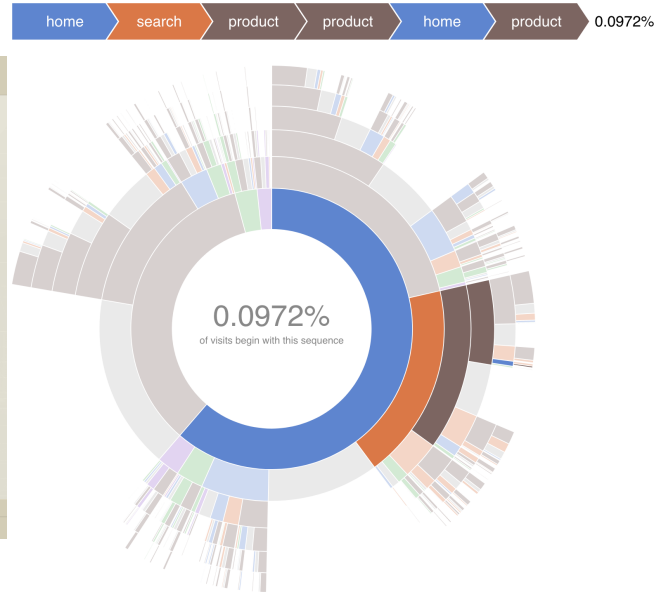
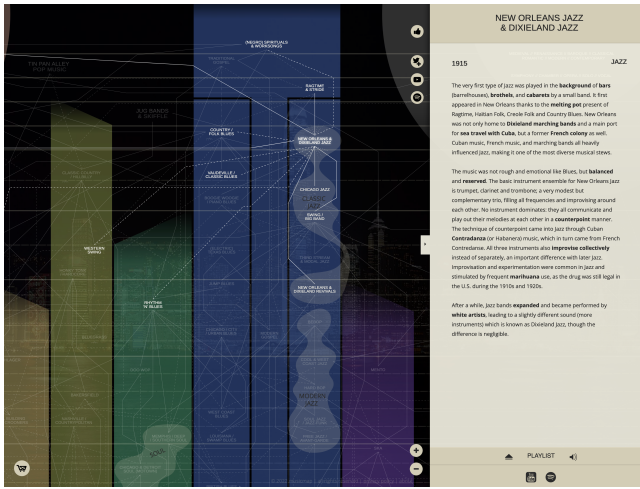


Figure 9.21: Animated and interactive charts II: The Genealogy and History of Popular Music Genres [↗](#), Musicmap, 2016 (top left); Sequences Sunburst [↗](#), Kerry Rodden, 2015 (top right); Health and Wealth of Nations [↗](#), Gapminder Foundation (middle left); Small Arms and Ammunition – Imports and Exports [↗](#), Google, 2012 (middle right); Möbius Transformations Revealed [↗](#), D.N. Arnold, J. Rogness, 2007 (bottom).

9.4 Miscellanea and Charts to Avoid

Some data visualizations are sufficiently unique that they cannot easily be grouped or categorized.

Chernoff Faces

Consider, as a singular example, **Chernoff faces**, which were designed on the premise that people can easily understand facial expressions. The Chernoff visualization can accommodate up to 18 or 36 facial feature variables.²⁹

29: The idea is perhaps intriguing and might even work well in some instances, but in most cases it fails to provide a useful rendering; among other issues, most facial features are not ordinal, faces are more than the sum of their parts, and not all facial features carry emotions.

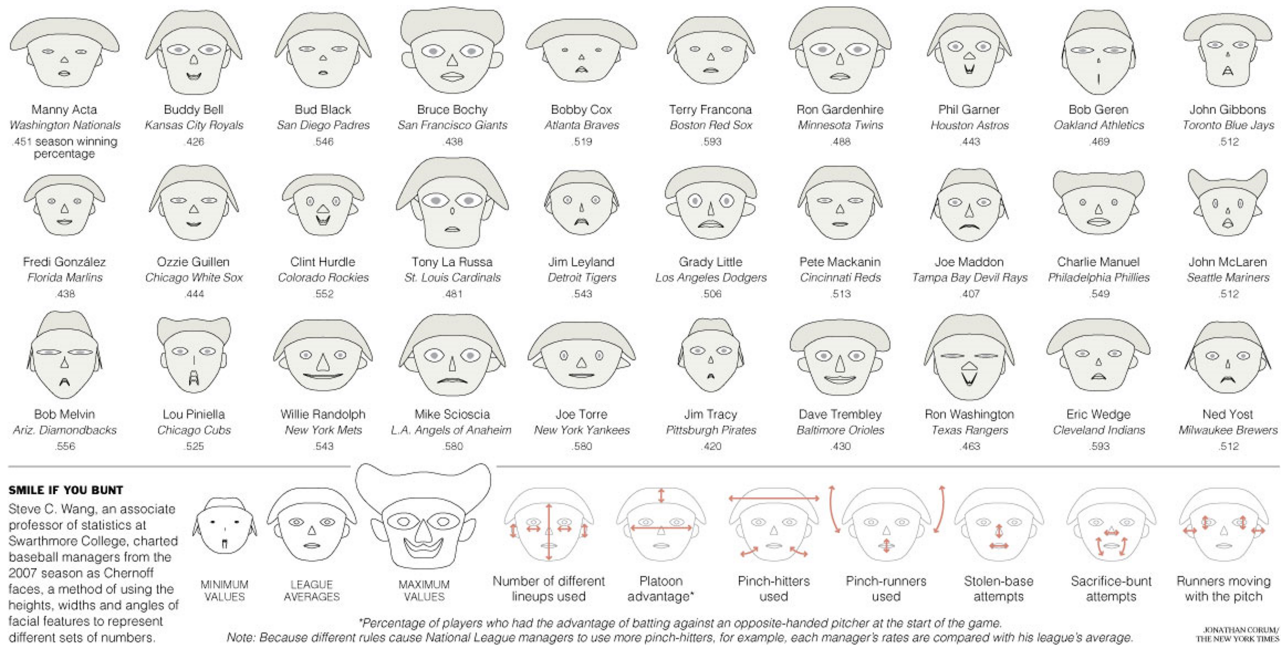


Figure 9.22: Chernoff faces of MLB managers characteristics during the 2007 season (SC. Wang, NY Times).

Alluvial and Sankey Diagrams

Alluvial and Sankey diagrams (see [here](#) and [here](#) for examples, respectively) are similar in appearance to one another, and both allow for the visualization of proportions; however, in the case of the alluvial diagram, the focus is on datasets with **multiple categorical variables**, and the chart displays the percentages of each variable relative to other variables.

Sankey diagrams, conversely, focus on **quantity breakdowns relative to particular categories** and how those quantities change when considering other categories.

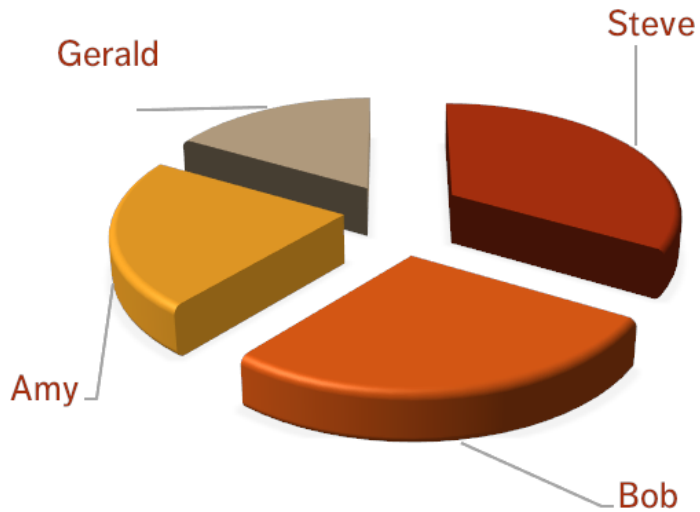
Charts to Avoid

One the one hand, we are agnostic when it comes to tools and methods: anything that helps convey the data story is on the table. On the other hand, some of the commonly-used approaches really put a damper on comprehension.

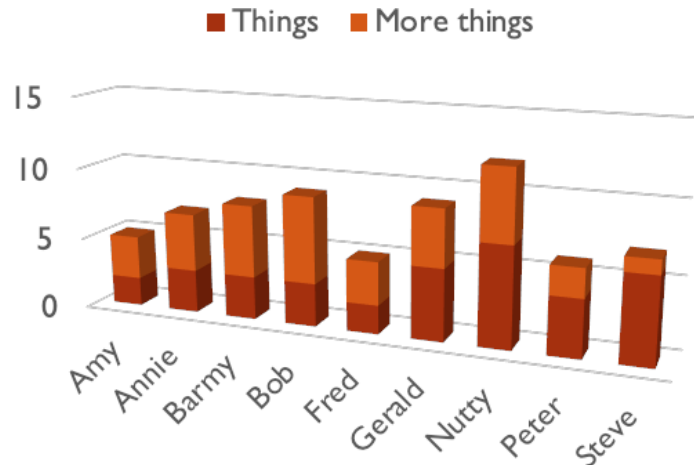
We strongly suggest avoiding:

- ANYTHING with an arc (except for gauge charts) such as pie and doughnut charts.³⁰ Human brains cannot easily compare angles and arcs, so these can become misleading: without labels, how easy is it to compare Steve & Bob below?

30: Sometimes we need to be pragmatic... but there are limits.



- **3D visualizations**, which we suspect are flat-out EVIL! As with arcs, we cannot easily visually compare data series in a 3D context; such charts are usually way too cluttered.



Note that there is always a danger that if certain types of visualization techniques take over, the kinds of questions that are particularly well-suited to providing data for these techniques will come to dominate the landscape, which will then affect data collection techniques, data availability, future interest, and so forth.