

Testing Statistical Charts: What makes a good graph?

Susan Vanderplas¹, Dianne Cook², and Heike Hofmann³

¹Statistics Department, Iowa State University, Ames, Iowa, United States, 50011;
email: srvander@iastate.edu

²Econometrics and Business Statistics, Monash University, Melbourne, Clayton
VIC 3800, Australia

³Statistics Department, Iowa State University, Ames, Iowa, United States, 50011

Xxxx. Xxx. Xxx. Xxx. XXXX. AA:1–29

[https://doi.org/10.1146/\(\(please add article doi\)\)](https://doi.org/10.1146/((please add article doi)))

Copyright © YYYY by Annual Reviews.
All rights reserved

Keywords

graphics, visualization, user testing, visual inference, perception, chart design

Abstract

It has been approximately 100 years since the very first formal experimental evaluations of statistical charts were conducted. In that time, technological changes have impacted both our charts and our testing methods, resulting in a dizzying array of charts, many different taxonomies to classify graphics, and several different philosophical approaches to testing the efficacy of charts and graphs experimentally. Once rare, charts and graphical displays are now everywhere - but do they help us understand? In this paper we review the history of graphical testing across disciplines, discuss different direct approaches to testing graphics, and contrast direct tests with visual inference, which requires that the viewer determine both the question and the answer. Examining the last 100 years of graphical testing, we summarize best practices for creating effective graphics, acceptance of the results of graphical testing, and what the future holds for graphics and empirical testing of interactive statistical visualizations.

Contents

1. INTRODUCTION	2
1.1. Design of Statistical Graphics	3
1.2. Statistical Mapping Using a Grammar of Graphics	6
2. TESTING METHODS	7
2.1. Explicitly Structured Graphical Tests	7
2.2. Implicit Graphical Tests Using Visual Inference	13
3. CURRENT BEST GRAPHICAL PRACTICE	16
3.1. Cognitive Principles	16
3.2. Chart Design	19
4. OPEN QUESTIONS AND FUTURE RESEARCH	22

1. INTRODUCTION

Any survey of literature on statistical charts and graphs will be complicated by the fragmentation of the literature, which occurs because any discipline which uses charts to summarize information generally also has individuals who research and assess the utility of these graphics. A quick survey reveals publications in journals from computer science, psychology, marketing and business, economics, ergonomics and human factors, statistics, sociology, communication and rhetoric, engineering, instructional design, and education. A review of early surveys of graphical forms suggests that this problem has long plagued the study of graphics, as Funkhouser (1937) discusses the broad disciplines affected by statistical graphics and Kruskal (1977) addresses the difficulty of a comprehensive review of the relevant literature.

Historically, the development of graphs and charts has been linked to the development of coordinate systems (Fienberg 1979, Beniger & Robyn 1978, Funkhouser 1937) and abstract representations of data. Preceding the development of formal mathematical coordinates, however, humanity has been representing information in abstract visual form since the origins of civilisation, e.g. spatial information using maps (Smith 1996) are displayed with varying degrees of abstraction. Figure 1 shows the oldest known world map, which dates from 6th century BCE Babylon.



Figure 1: *Imago Mundi* Babylonian map, which is the oldest known world map (6th century BCE). The representation of the world is relatively abstract. The world is shown as a disc, with Babylon represented by a rectangle at the top end of the Euphrates river, which flows south to the border of the disc. Several other population centers are marked with small circles. (The British Museum 1882)

Image courtesy of the British Museum, released under a CC By-NC-SA 4.0 license.

Exploring visual abstractions of data, has been an active pursuit of the “polymaths” who, grappling with increasing collections of economic, demographics, and measurements on the natural world, developed the foundations of current scientific pursuits. During the 18th and 19th centuries, governments became involved, assembling data and graphical representations into the statistical atlases, and government issued reports, as created by Harms (1991), Playfair (1801), Walker (1874), and many other contemporary workman. In the 20th century, corporations began using charts and graphs to understand their inner workings - studies of the use of charts and graphs at AT&T (Chandar et al. 2012) and DuPont (Yates 1985) show efforts to standardize and formalize the use of graphics in decision-making at both companies.

As new charts were invented to represent data differently and highlight features of data (Bernstein & Cowden 1937, Yates 1985, McDonald 2014), discussions about the use of statistical graphics began to appear in the literature (Peuchet & Gilbert 1805, Brinton 1917, Karsten 1923), including the relative strengths and weaknesses of various types of charts. In most cases, the drive to produce a classification system for charts and graphs or a system of recommendations for presenting charts and graphs were based on heuristics and largely unsupported by experimentation (Kruskal 1977, MacDonald-Ross 1977). Many of these ad-hoc classification systems could not accommodate the large numbers of new plot types being developed.

Calls for experimental validation of the perception and utility of statistical charts were heeded, though at first the experiments were fraught with methodological issues (Croxton & Stryker 1927). Much of the early experimentation regarding the accuracy of graphical forms was based in psychophysics research (Teghtsoonian 1965) on the perception of size and shape. Eventually experiments became more naturalistic: cognitive psychologists and statisticians began testing different types of graphics, identifying types of perceptual errors associated with different plots (Spence 1990, Cleveland & McGill 1985). In most cases, this testing was limited to simply reading information from the charts, using accuracy or response time measurement. More recently, other methods for examining statistical charts have been developed, including the lineup protocol (Wickham et al. 2010). Even with these developments, the aim of most experimental research in statistical graphics focuses on the initial perception and graph comprehension. Very little work has been done to understand the effect of charts and graphs on higher cognitive processes such as learning or analysis (Green & Fisher 2011).

All the graphic methods enumerated already exist and it is unlikely that a completely new one will be invented.

Since they have been used separately up until the present time it would be convenient to classify them by analogy, according to the categories above... -

International Statistical Congress (1858), in Funkhouser (1937)

1.1. Design of Statistical Graphics

Charts and graphs are used for many purposes (Tukey 1972, Fienberg 1979): to summarize data, for analysis, exploration and discovery, diagnosis of statistical relationships, to make a rhetorical argument, or even as a substitute for tables. The initial heyday of graphic design was enabled by colour lithography used charts and graphs to tell stories about nations and events (Kostelnick 2016), but in the first half of the 20th century, graphics were regularly used for mundane purposes as well, such as supporting business decisions (Chandar et al. 2012, Yates 1985) and communicating weather forecasts. As technology has developed, allowing charts to be created quickly for exploratory purposes, the gaps between graphics for presentation, entertainment, and analysis have widened. The different purposes motivating the creation of the chart influence the form and complexity of the chart, and the intended audience and the reach of the chart are also important considerations.

It is useful to consider a continuum from utilitarianism to artistry, where purely utilitarian charts are, as advocated by Tufte (1991), devoid of “chartjunk” or any decoration, and purely artistic charts trade accurate representation of the data for visually compelling renderings. The distinction between infographics and statistical graphics is primarily one of intent: the infographic, often composed of many small simple plots interspersed with pictures and text, is designed to attract attention and tell a story. In contrast, the statistical chart is designed to effectively and accurately show the data, potentially with accompanying statistical model information – any visual enhancements should contribute to that aim (Gelman & Unwin 2013, Wickham 2013).

Figures 2, 3 and 4 show plots designed for different statistical purposes. The Hertzsprung-Russell diagram allowed astrophysicists to make a convincing argument about the lifecycle of a star, based on its temperature, colour, and size. The hurricane forecast map is designed to communicate urgent information and facilitate decision-making by the US National Hurricane Center. It shows forecasts of the hurricane’s path, the wind field, forecasted strength, and coastal warnings. Figure 3 shows the prediction map for Hurricane Michael, which hit the Florida panhandle in October 2018. It is primarily utilitarian, with additional decorations such as ocean colour and state boundaries provided for geographic context. The birthday chart reproduced in Figure 4 shows the average number of births in the United States on each day of the year, organized by month. Figure 4 is a static reproduction of an interactive plot used as part of a web page (<http://thedataviz.com/2016/09/17/how-common-is-your-birthday-dailyviz/>) telling a story about birth dates. It is simple in form, almost infographic style. By engaging with the interactive chart, a reader could easily deduce that there are relatively few babies born on major holidays, such as Memorial Day, Independence Day, and Christmas; there are also considerably more babies born in the summer months than in the winter.

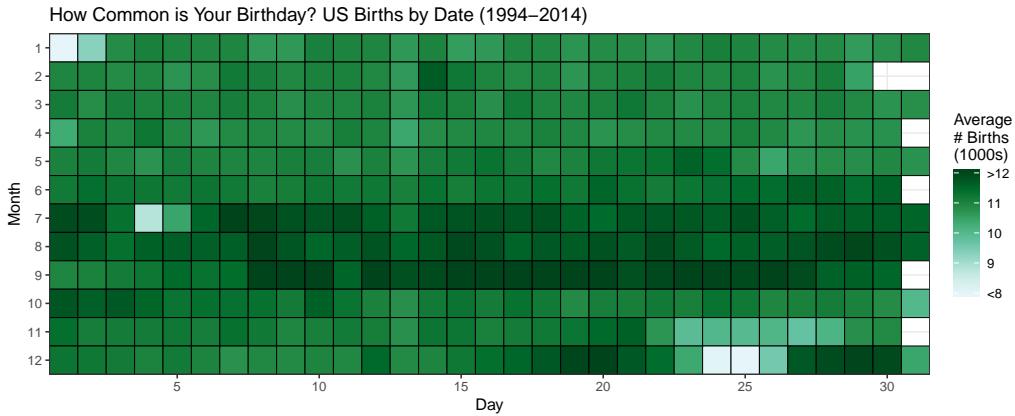


Figure 4: A static reproduction of the birthday chart, which shows the average number of births on a yearly basis in the United States, by month and day of the month. The interactive version, available at [http://thedataviz.com/2016/09/17/how-common-is-your-birthday-dailyviz/](http://thedataviz.com/2016/09/17-how-common-is-your-birthday-dailyviz/), is designed to be a storyteller to encourage the viewer to interact, check the estimated conception dates and read the frequency. Data source: CDC National Center for Health Statistics (1994-2003) and Social Security Administration (2000-2014), as reproduced at <https://github.com/fivethirtyeight/data/tree/master/births>.

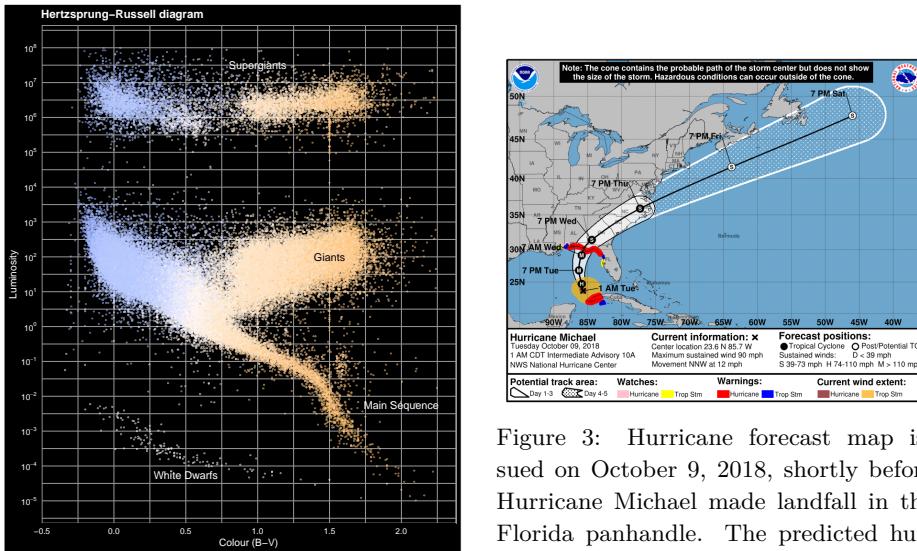


Figure 2: The Hertzsprung-Russell diagram (circa 1910) allowed astrophysicists to make the connection between a star's temperature, visual colour, and size, enabling a better understanding of the life cycle of a star.

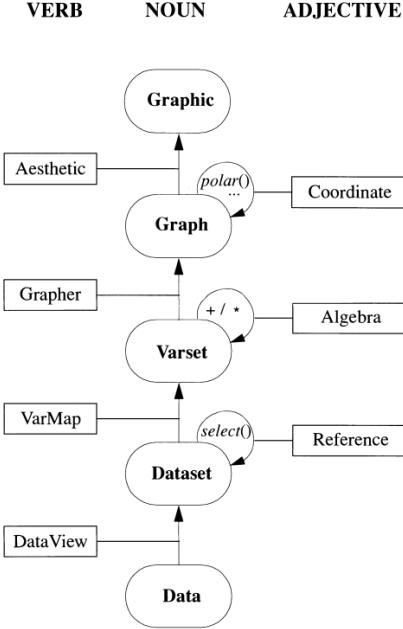


Figure 5: From data to graphic, from (Wilkinson 1999, Fig 2.1). The schematic representation of the steps required to create a graphic from a dataset, which can be used to specify the variable mapping, data transformations, coordinate system, and aesthetic features independently.

1.2. Statistical Mapping Using a Grammar of Graphics

In parallel with efforts to understand the perception of charts, there have been many attempts to develop systems for classifying graphics, including Bertin & Berg (1983), Desnoyers (2011), and Wilkinson (1999). Systems which attempt to categorize charts based on their geometric representations generally make no effort to include all types of graphics, and have difficulty accommodating charts which may fall into two or more categories. The classification of graphics based on the underlying components and their relationships, as in the grammar of graphics developed by Wilkinson are more robust; they also provide an elegant framework for comparing different types of graphical representations separate from the underlying data structure. An analogy to conceptualise the difference is that the former is like treating plots like creatures in a zoo, with a unique name for each, while the latter is analogous to having a phylogeny based on genetic data showing how plots are related.

Figure 5 shows the framework of the grammar of graphics, where the data is filtered, variables are mapped, transformations are specified, and then finally, transformed data are mapped to plot aesthetics and coordinate system specifications to produce an abstract visual representation of the data. Full or partial implementations of the grammar of graphics are available for most common scientific computing languages: e.g. `ggplot2` in R (Wickham 2010), `plotnine`, a python implementation of `ggplot2` (Kibirige 2017), and `Gramm` in Matlab (Morel 2018).

The grammar of graphics also enables data plots to be considered to be statistics (Ma-

jumder et al. 2013). A statistic is a functional mapping of a variable or set of variables. With “tidy data”, that is, data where each variable is in its own column, each observation is in its own row, and each value is in its own cell (Wickham & Grolemund 2017), the grammar of graphics creates visual statistics. Variables, as columns in the data table, are mapped to graphical elements, such as the x axis, or y axis, or to colour, shape or even facet, using the grammar. The data plot can then be treated like other statistics: by imagining what the plot might look like in the absence of any structure, we can use the plot of the actual observed data to test for the likelihood of any perceived structure being significant.

Using the grammar of graphics, it is easy for experimenters to compare different types of charts using the same data, as the underlying structure of the graph remains the same. Figure 6 shows three plots created using the same data and different geometric objects, with the `ggplot2` code to create the plots. Comparing these graphics experimentally would be reasonably simple as the grammar of graphics helps to control the extraneous variables introduced by utilizing different plot types. In addition, this approach to transformations and scales allows experimenters to easily test judgments made utilizing different axis transformations and colour scales to compare perceptual accuracy (Hofmann et al. 2012, VanderPlas & Hofmann 2017).

2. TESTING METHODS

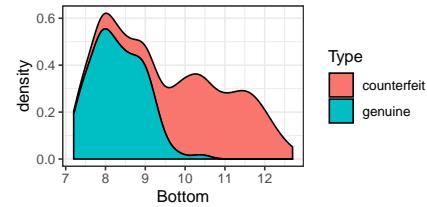
In this section, we will distinguish between explicitly structured graphical tests, which require the participants to answer specific questions about the graphical objects under experimentation, and implicitly structured tests, where the participant must infer the questions of interest from the provided stimuli.

2.1. Explicitly Structured Graphical Tests

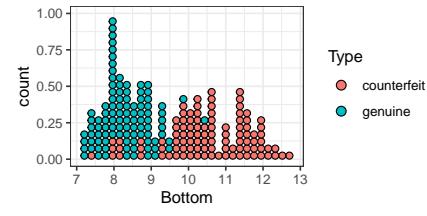
During the perceptual process, information from a visual scene is processed by the brain, with information extracted at many different levels of the cognitive process. Preattentive perceptual effects are those which do not require sustained cognitive attention; they are processed automatically within the first 500 milliseconds of viewing a chart or graph. Components processed preattentively include colour and shape, as well as some basic information about coarse relationships between individual components. After the preattentive stage, attention is necessary for subsequent processing; this directed attention scaffolds relationships between components and helps us interpret the chart or graph in context. Most of the insights we gain from charts and graphs are due to the cognitive processes that occur after attention is focused on specific aspects of the graph; as a result, most of the testing methods we will discuss are focused on the attentive portion of the perceptual process.

2.1.1. Preattentive Graph Perception. Initial research into preattentive perception used a search task, where participants had to identify a particular object in a field of distractors, manipulating display size and varying one or more features such as colour or shape; participants’ search times were measured to determine the amount of effort necessary to complete the search task. Preattentively perceived features showed a near constant reaction time over increasing display size, while features which are processed attentively show an increasing reaction time with increased display size (Treisman 1980). A primary question in the discussion of preattentive graph perception is whether there are advantages in designing a

```
# Density plot
ggplot(data = banknote,
       aes(x = Bottom,
            group = Type,
            fill = Type)) +
  geom_density(position = "stack")
```



```
# Dotplot
ggplot(data = banknote,
       aes(x = Bottom,
            group = Type,
            fill = Type)) +
  geom_dotplot(method = 'histodot',
               stackgroups = TRUE)
```



```
# Histogram
ggplot(data = banknote,
       aes(x = Bottom,
            group = Type,
            fill = Type)) +
  geom_histogram(position = "stack")
```

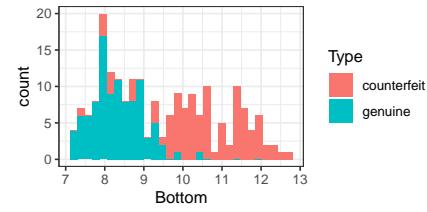


Figure 6: Three different plots of the Swiss banknote data (Flury & Riedwyl 1988), created using the grammar of graphics as implemented in `ggplot2`. The data consist of measurements of the dimensions of the banknotes, including the length of the bottom edge, which is shown in the plots. The main plot specification and syntax remain the same, with the form of the plot changing due to the specification of the geometric object used to represent the data.

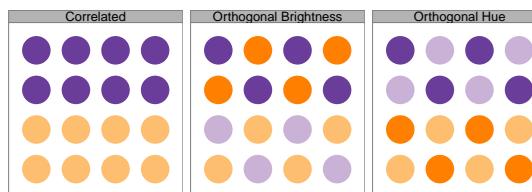


Figure 7: Callaghan (1984) showed irrelevant variation increases response times, indicating that hue and brightness are integrated and preattentively perceived. The images shown here are examples of the types of stimuli used in the experiment. While hue and brightness can be manipulated separately, they are preattentively perceived as a single unit.

graph to promote the preattentive perception of features, ideally, reducing cognitive load. We will distinguish between tests which use graphical forms and more primitive tests which

use basic geometric elements during the testing process. The results from more primitive experimental designs still apply to the design of graphs and charts, but the experimental design does not involve any display of actual data (Figure 7 is an example of an experimental stimuli that does not represent actual data). Preattentively processed features include shape, angle, size, and texture; however, typically, combinations of preattentive features which represent separate features in the data are processed attentively, with at least one major exception. Callaghan (1984) demonstrated that hue and brightness are integrated, that is, that even though they can be separately manipulated, they are still perceived preattentively as a single unit, using arrays of tiles similar to those shown in Figure 7.

Extending this work, Healey et al. (1996) has used the same segmentation paradigm, applied to more complex charts utilizing actual data, with the goal of exploring region segmentation using preattentive cues. Healey's experiments use multivariate displays, leveraging the preattentive grouping of similar objects to separately represent features using colour, height, and texture (Healey & Enns 1999). While these displays may not follow best graphical practice in other respects, they do show the utility of designing with the preattentive perceptual process in mind.

2.1.2. Attention Mediated Testing Methods. Creators of a chart or graph typically operate under the (hopefully safe) assumption that readers will spend more than 300 milliseconds considering its contents; as a result, attention mediated testing methods allow a more realistic mechanism for testing overall performance of different graphical forms. Cleveland & McGill (1987) discuss the different approaches to research methodology in this area, and while they do not include all of the experimental approaches we discuss, the paper makes an important distinction between informal and formal graphical exploration. In the informal approach, changes are made to the graph and the iterative versions are compared to determine what information is easily accessible; in the formal approach, an experiment is designed and participants are tested in a controlled manner. This section describes several different experimental approaches which can be used to answer the general question of “how effective is this graph at communicating useful information?”

2.1.2.1. Direct Observation: Numerical Estimation, Speed, and Error Rates. One of the simplest ways to test the utility of a graph is to verify that information can be accurately read from it. Charts are, after all, generally recognized as having more utility than tables for presenting information in an accessible and useful way; if that information cannot be read back out in a relatively accurate manner, the graph's utility is suspect. Early experiments, such as Eells (1926), Croxton & Stryker (1927) and Croxton (1932) used accuracy alongside speed and other considerations for plot evaluation. Later studies (Peterson & Schramm 1954, Cleveland & McGill 1984, Broersma & Molenaar 1985, Dunn 1988, Tan 1994, Amer 2005) were conducted with similar methodology; in essence, the participants are provided with a chart and asked to estimate some quantity or answer a predefined question using the information provided in the chart. The accuracy of various types of charts, as measured by participants responses to the questions, is then used to determine which charts are superior. It is important to ensure that the specific charts and questions used are aligned; studies are commonly critiqued on the basis that the charts or the questions were not appropriate for the task — indeed, the first studies of pie vs. bar charts were heavily criticized on these grounds (von Huhn 1927).

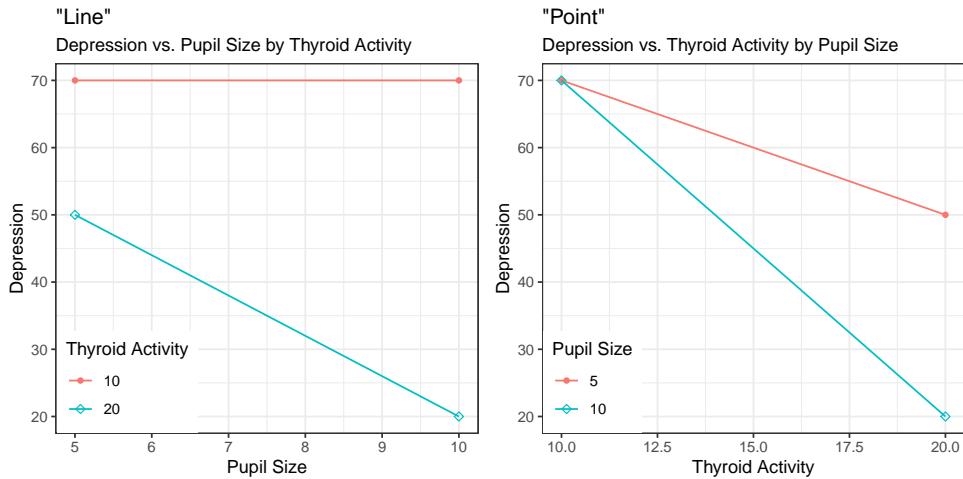


Figure 9: Sample plots from Shah & Carpenter (1995), showing data where one independent variable (pupil size) has no effect on the y variable (depression) for one value of the second independent variable (thyroid activity). The two representations of this data, referred to as line and point, are not equally effective for communicating the joint relationship of the two independent variables.

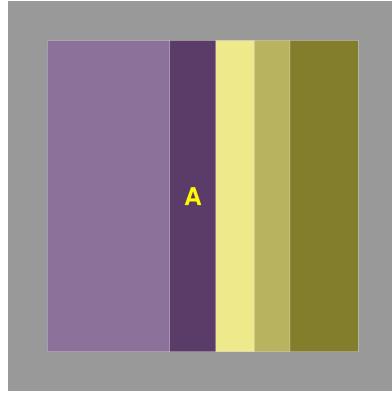


Figure 8: A framed spine plot similar to those shown in the 1874 Statistical Atlas, recreated using modern methods. The chart shows the proportion of church accommodations in Oregon, with the four largest denominations in the state shown in the center and the unaccommodated population shown in the grey region framing the plot. Participants were asked to estimate the size of the area labeled “A” relative to the entire chart; the presence of the frame significantly decreases the accuracy of these estimates. (VanderPlas et al. 2019)

A similar type of study avoids the pitfalls associated with numerical estimation by showing participants multiple charts in sequence, and asking them to evaluate the differences in the charts with respect to the dimension of interest. Shah & Carpenter (1995) required participants to examine the relationship between three variables as shown in several different types of line charts. One set of charts used in the experiment is recreated in Figure 9; in the first plot, it is much easier to determine that there is no relationship between pupil size and depression for one level of thyroid activity compared with the second plot.

In some studies, response time is the main quantity of interest. Participants are provided with a certain chart or stimulus, and the amount of time spent considering the stimulus before answering are used together with accuracy of the answer to assess the difficulty of the task, under the premise that more difficult or mentally demanding tasks will require more

time before a response is generated. Carswell & Wickens (1987) uses the error rate and response time to determine whether triangular “object displays” are superior to bar graphs to represent the inputs and outputs of a dynamic system (see Fig.2 in Carswell & Wickens 1987, for an example of an object display); Legge et al. (1989) uses efficiency, a function of time and accuracy, to assess different types of charts. More modern studies may collect both response time and numerical estimation data online, using services such as Amazon Mechanical Turk to conduct usability studies in statistical graphics (Heer & Bostock 2010). We have even used direct estimation in combination with online testing to examine very old charts from the 1874 Statistical Atlas (VanderPlas et al. 2019), determining that framed mosaic, spine, and pie charts are sub-optimal graphical representations. Figure 8 direct estimation used in the experiment; participants entered an estimate of the proportion of the chart represented by the area labeled “A”. It should be noted that some of the other experimental paradigms, such as psychophysics and implicit tests, can also be used in combination with Amazon Turk and other online testing services.

There are limits to what one can test using direct estimation: it is generally preferable to test only very straightforward assessments of the content of a chart or graph, to fit within a simple experimental paradigm. In addition, open-ended estimation tasks elicit certain well-known biases such as the tendency to round to multiples of 5 or 10 (Baird et al. 1970). Long-term interaction with a complex graph or chart showing multiple layers of data is generally not ideal within this paradigm, which requires a fixed set of numerical assessments that do not accurately represent how we explore a new, complex graphic. This complexity explain why there are so few studies of rich, complex graphics that may require domain expertise in addition to the ability to read information from a visual display. To approach situations with more complicated graphics, or charts which are known to induce perceptual biases in the participants, it is often more useful to utilize other experimental paradigms that facilitate examination of specific parts of the perceptual process - such as the use of eye tracking to measure attention and motivation, or the use of verbal descriptions to assess more complicated graphs. The next sections cover some of these more nuanced approaches to explicit testing of graphs.

2.1.2.2. Psychophysics and Signal Detection Theory. Some studies of graphs utilize psychophysics methodology to assess data visualizations. Initially, of course, a significant portion of the research in statistical graphics came from the fields of psychophysics and cognitive psychology (Spence 1990, Teghtsoonian 1965, Lewandowsky & Spence 1989), but in most cases this was not accompanied by a use of the methods of psychophysics for experimental testing of charts and graphs. Psychophysical experimental design is focused on whether an effect is detectable, and whether the magnitude of the effect can be accurately estimated. Common methods, such as the method of constant stimuli and the method of adjustment, involve repeatedly presenting a participant with charts and asking them to evaluate the chart on the basis of a particular question of interest. In the method of adjustment, this is done with the control of the participant, who adjusts the stimuli interactively until the effect is just barely noticeable; in the method of constant stimuli, the effect size changes randomly from trial to trial to reduce continuity effects. In graphical testing, Hughes (2001) used the method of constant stimuli to assess the ability to detect a difference in height in 2D or 3D barcharts; VanderPlas & Hofmann (2015) used the method of adjustment to experimentally determine the size of the line-width illusion’s distortion of variance perception. Psychophysics methods also seem to be relatively common in studies

of map perception, particularly when the goal is to estimate the amount of exaggeration or other corrective distortion necessary for realistic perception of the map (for an overview of relevant cartography studies, see Brandes 1976).

2.1.2.3. Thinking Aloud. Another approach to testing graphics is to examine the cognitive processes which occur as a graph is read. Lacking mind-reading devices, the next best option is to ask participants to talk through their thoughts as they read and use a graph in a realistic setting (Concurrent Think Aloud, CTA), or as they recall a graph after the fact (Retrospective Think Aloud, RTA) (Guan et al. 2006). The think aloud process allows experimenters to examine the use of complex graphics “in the wild” or at least in situations that are less artificial than the paradigms allowed by numerical estimation and psychophysics methods. Think-aloud studies allow researchers to attempt to measure insight (North 2006) and reasoning (Dunbar 1995) in complex situations such as experimental design, decision making (Normand & Bailey 2006), or the process of weather forecasting (Trafton et al. 2000, Kirschenbaum 2003). In forensics, think-aloud studies are known as “white box” studies, because it is possible to “see” what a forensic examiner is thinking and why they make a specific conclusion about the evidence (Ulery et al. 2011). Studies using think-aloud protocols have also examined the process of exploratory data analysis, finding that unexpected results are more likely to be represented in informal terms initially, but that with familiarity language shifts to formal explanations (Trickett et al. 2000a,b). The think-aloud protocol is also conducive to use with interactive graphics, combined with logging or video recording software which can record the state of the graphics device in parallel with the user’s monologue. Studies have also combined think aloud protocols and eye tracking studies with the goal of validating CTA (Cooke 2010) and RTA protocols (Guan et al. 2006) for use with statistical graphics and general usability testing. While the data which results from the think-aloud protocol is typically more qualitative and less quantitative than results produced using other methods, there is significant additional insight into the underlying cognitive processes affecting visualization which cannot be obtained through other means.

2.1.2.4. Eye Tracking. Where think-aloud protocols allow insight into the cognitive process, eye tracking facilitates insight into the process of visual attention, providing data on the approximate spatial location of visual focus. The attention-fixation process occurs too quickly to be accurately verbally communicated, but eye-tracking equipment allows experimenters to identify the portions of the chart which require attention and sustained cognitive effort or which attract interest from participants, inferring from gaze and fixation the cognitive processes occurring during graph comprehension. Eye tracking allows researchers to determine that viewers spend relatively little time examining the axes in scatterplots, but significant amounts of time examining the axes in parallel coordinates plots (Netzel et al. 2017), suggesting that the process of reading these two chart types is fundamentally different. Another study leveraged eye tracking to identify features which provide useful information during the graph reading process for several different types of charts (Goldberg & Helfman 2010). Eye tracking is a powerful tool when combined with good experimental design: Fabrikant et al. (2010) examined fixations when shown meteorological charts before and after users are provided with introductory training about meteorology, finding that after training users conclusions were more accurate, response time increased, and fixations were directed to more useful areas of the maps. The use of eye tracking with different

source populations also allows researchers to understand how dyslexia (Kim et al. 2014) and graph literacy (Woller-Carter et al. 2012, Okan et al. 2016) affect the graph comprehension process, providing better design guidelines for specific target audiences.

2.1.2.5. Combination Experiments. Many studies use a combination of the explicitly structured graphical tests discussed here. Think-aloud studies are relatively easy to integrate with eye tracking studies, and it is not difficult to add in direct observations or psychophysical evaluation methods as part of the trial design. Psychophysics and direct observation studies are limited by the questions that are asked; eye tracking results only provide information on visual focus; think-aloud results are generally qualitative, but when these techniques are combined, they provide a much more complete picture of the cognitive processes which underlie graphical perception. Ryu et al. (2003) used a combination of eye tracking, cognitive tests, direct observation, and think-aloud protocols to examine the integration of information across multiple types of charts, determining that integrating information from multiple parallel coordinates plots is slow, difficult, and inaccurate compared to information integration when a scatterplot or map is presented with a parallel coordinate plot. Zgraggen et al. (2018) also used a combination of think aloud, eye tracking, direct observation, and interactive graphics to examine the impact of exploratory data analysis on multiple testing problems. Many of the attention-mediated, explicitly structured testing methods can also be combined with implicitly structured tests, which we discuss in the next section, to produce a more comprehensive view of the process of graphical perception.

2.2. Implicit Graphical Tests Using Visual Inference

Explicit graphical tests, as we have referred to them, are tests where the user is directed to assess a specific feature of a plot or answer a specific question. That is, the tested hypothesis is explicitly stated, providing the user with cues to the intended purpose and function of the plot and/or the relevant features of the data shown within the plot. In contrast, in an implicit graphical test, the user must identify both the purpose and function of the plot and use that information to evaluate the plots as shown. Typically, these tests are structured as visual inference problems, as introduced in Buja et al. (2009), though other formulations of implicit tests exist as well (Hasanhodzic et al. 2010). Explicit tests are typically conducted on plots which have been created to showcase specific structure in the data in order to present results; in contrast, implicit tests are designed to inform exploratory data analysis (as advocated by Tukey 1965) and the iterative model diagnostic process. During a data exploration (EDA), statisticians typically examine the data using many different plots considering different aspects of the data and selecting only those results which are visually interesting for further exploration. This raises the question whether the ‘visually interesting’ features in those plots show actual signal and do not arise simply by chance. This question is a direct consequence of our general use of graphics and our definition of ‘interesting’. Generally speaking, something is ‘interesting’ that does not follow our expectations in some way – in other words, we conduct an implicit hypothesis test when drawing and looking at charts; ‘interesting’ charts are those which are deemed ‘significant’ according to this implicit hypothesis test. Making this hypothesis test explicit allows for a more formal evaluation of the significance of a visual finding. In terms of a classical hypothesis test, a plot of the data is taking the role of a (visual) test statistic. Null plots are created from data in accordance with the null hypothesis (e.g permuted data, if the null

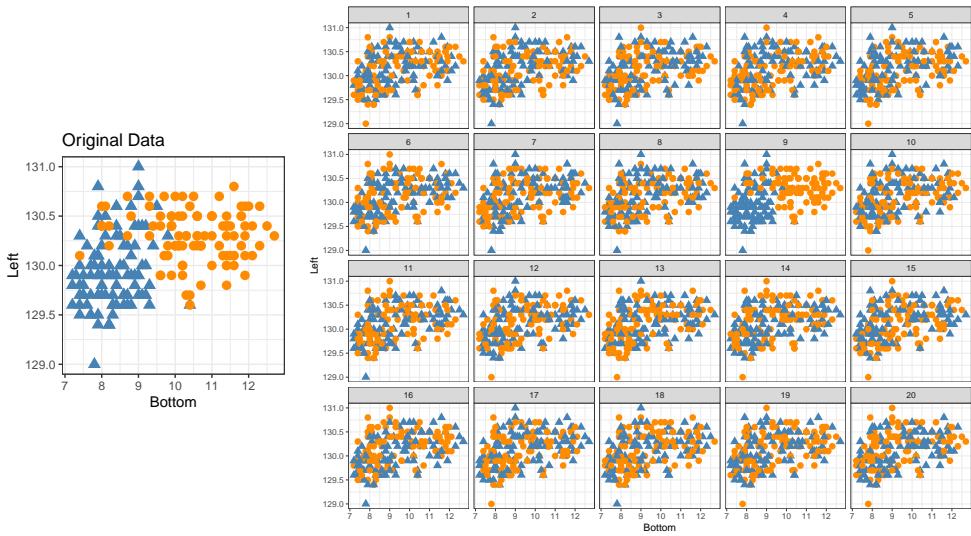


Figure 10: The original data (left) is embedded into a set of 19 generated under the null hypothesis. A visual hypothesis test is conducted by asking multiple users to select one (or more) plots which are different; a p-value can then be calculated from the answers. In this particular lineup, plot 11 contains the original data; the other plots are generated by randomizing the vector of colour/shape values as in a randomization test.

hypothesis assumes that there is no non-random relationship between the variables). If the data plot is visibly different from the null plots, i.e. can be picked out through visual inspection by an independent evaluator, this counts as evidence against the null hypothesis and once enough evidence is accumulated we would conclude that the null hypothesis has to be rejected and accept that the feature in the data plot is not random. Formally, a statistical lineup of m visual statistics consists of $m - 1$ plots $T(y_0)$ simulated from the model specified by the null hypothesis, and the test statistic $T(y)$ produced by plotting the actual data, which may arise from the alternative hypothesis. Figure 10 demonstrates this process, with the 19 null plots generated by randomizing the values of y . The statistical lineup is then evaluated by K independent observers, with the resultant p-value calculated according to the null hypothesis that each of the m visual statistics are equally likely to be selected (Majumder et al. 2013). The `nullabor` package contains tools for graphical inference (Wickham et al. 2018), including p-value calculations and power calculations for visual inference.

Implicit graphical tests approach the problem of spurious plot relationships at the level of the data, leveraging the human visual system to conduct a suite of visual tests for features such as outliers, clusters, linear and nonlinear relationships. The advantage to implicit testing is that lineups do not require a specification of a feature of interest in the testing framework, i.e. we do not have to ask questions such as “which group has a higher proportion of responses”. Much of the historical research of comparing different types of charts has been criticized because the specific question phrasing does not provide readily generalizable results; the lineup protocol removes this obstacle by charging the user with the task of identifying the most different looking plot and thereby selecting the feature with

the visually most salient difference compared to the other plots. This allows an evaluation of competing plot designs without the complications of potentially steering participants towards a particular outcome by phrasing questions. This leads to a completely data-driven result: if bar charts are indeed better suited for a task than pie charts, the target plot will be selected more frequently when the lineup is presented with bar charts than when the lineup is presented with pie charts. The real power of the lineup protocol is that when *combined with the grammar of graphics*, we can hold the underlying data and summary statistics constant, isolating the effect of different plot types, coordinate transformations, and aesthetic mappings on our ability to detect effects in the data.

Statistical lineups have been used experimentally to examine single plot types in many contexts: residual plots in hierarchical linear models (Loy & Hofmann 2015), perceived clustering in high-dimensional data (Roy Chowdhury et al. 2015), and spatial clustering in geographical research (Widen et al. 2016). Loy et al. (2016) used statistical lineups to evaluate different types of Q-Q plots for assessing violations of normality, determining that the visual tests were generally more powerful than common numerical tests when assessing violations of normality. Beecham et al. (2017) used lineups to assess the effect of spatial autocorrelation when represented using different grid structures, finding that lineups can be used for chloropleth maps if the null plots are generated under models with reasonable spatial autocorrelation models. Other studies have also found that the approach to null model generation is critical. It can be difficult to specify the null data generating model in a way that adequately mimics the data plots, which suggests that visually we are able to identify many more features than those typically tested using standard quantitative hypothesis tests. This implicit testing of many different hypotheses does make null distribution specification challenging, but also highlights the power of visual cognition to detect subtle differences in data.

Typical lineups contain a single target plot, but this is not a requirement. VanderPlas & Hofmann (2017) used two targets, each generated from a competing data model: clustering, or linear association. The null plots in this experiment were composed of a mixture of points generated from each of the two data models, ensuring that the two targets were both slightly more extreme than the null plots, but that the null plots and the target plots shared some features. Using this approach, the authors tested the effect of different aesthetics on the selection of each of the two targets, examining the strength of aesthetics such as colour, shape, trend lines, error bands, and 95% ellipses for highlighting clustering or linear trends in the data. By providing viewers with a choice between data generated by competing models, the two-target lineup approach provides a way to directly examine the visual strength of each model compared to the null and comparing the models directly. For instance, when comparing a model generating clustered data to a model generating data with a linear association between y and x , this protocol establishes that colour and shape aesthetics slightly increase the likelihood that the cluster target plot is selected, while a trend line aesthetic slightly increases the probability that the linear relationship target plot will be selected. More broadly, VanderPlas & Hofmann (2017) shows that the aesthetics used in a plot can significantly impact the perceived relationship between variables.

In the ten years since the introduction of the lineup protocol, many studies have leveraged the grammar of graphics to ensure that the underlying data mapping remains the same while manipulating the geometric representation of the data and overall visual appearance of the plot. Combining the conclusions from both implicit lineup tests and the explicit tests described in the previous section, what can we say about best graphical practices, beyond

“pie charts are awful”?

3. CURRENT BEST GRAPHICAL PRACTICE

All of the user testing in the world cannot identify the “best” possible graphic – we can instead only experimentally assess which graphical designs are better for a specific purpose. This can lead to a rather fragmented approach when describing “best practice”, and so in order to avoid this, we will examine graphical practice using the principle of “First, do no harm” from the Hippocratic oath.

3.1. Cognitive Principles

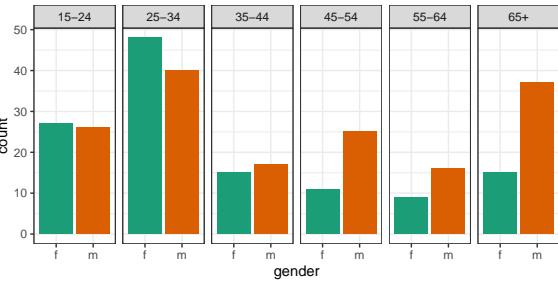
A useful starting point is to apply Gestalt principles of visual perception (Wagemans et al. 2012a,b), such as proximity, similarity, common region, common fate, continuity, and closure, to data plots. These principles are useful because good graphics take advantage of the human visual system’s ability to process large amounts of visual information with relatively little effort. Understanding the principles which underlie this processing allows us to create charts which require less cognitive effort to read, freeing us to think about the content, rather than the form of the chart.

3.1.1. Proximity. The principle is ‘objects or shapes that are close to one another appear to form groups’. For plot design, proximity is used to place items to compare close together, and less important comparisons further apart. Figure 11 illustrates this principle when using facetting in a plot, using data on tuberculosis (TB) incidence in Australia in 2012. When plots are *facetted by age*, it means gender is easier to compare. We learn that females more than males are detected to have TB in their early 20s, but in the aging population, males more commonly are detected with TB. Conversely, when plots are *facetted by gender*, the distribution of age is easier to examine. We learn that the age distribution of TB incidence for females skews heavily towards younger women. For males, TB incidence is more uniform across ages because there are high counts at young, old and middle age categories. Effectively utilizing proximity in organising plots, makes a huge difference in ease of information communication.

3.1.2. Similarity. The gestalt principle of similarity suggests that we group things which have similar appearance and exclude objects which have a different appearance. In charts, this principle is often leveraged by colouring points or bars according to a categorical variable, or by using points of different shapes to represent different categories. VanderPlas & Hofmann (2017) showed that the addition of colour and shape to a scatterplot increases the likelihood that individuals will perceive clustered groups of points. In Figure 11, the colouring of bars allows us to easily see that the similarly coloured rectangles represent the same group of people, even though the bars are separated by facets and other groups.

3.1.3. Common Region. The gestalt principle of common region suggests that elements contained within a common region belong together. Common region helps us to easily read small multiple plots, because the graphical elements of each small plot are grouped into a single entity that can be examined on its own. In addition, confidence bands and bounding ellipses also activate this gestalt principle by grouping points within the boundaries

```
# gender in age
ggplot(tb_au_12,
       aes(x = gender,
            y = count,
            fill = gender)) +
  geom_bar(stat = "identity") +
  facet_grid(~ age)
```



```
# age in gender
ggplot(tb_au_12,
       aes(x = age,
            y = count,
            fill = age)) +
  geom_bar(stat = "identity") +
  facet_grid(~ gender)
```

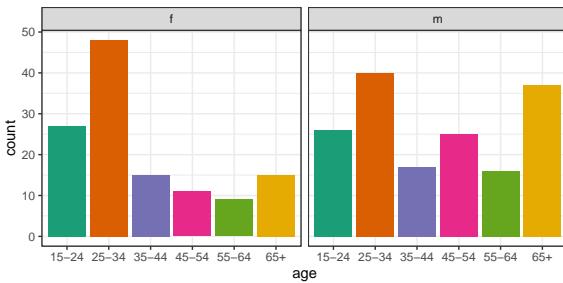


Figure 11: Two different arrangements of the same data, to illustrate how the proximity of elements makes a question easier to answer.

together (VanderPlas & Hofmann 2017), highlighting the presence of outliers which do not belong to the main group.

3.1.4. Common fate. The gestalt principle of common fate describes the tendency to group objects which are moving together in the same direction and at the same speed together. Common fate is certainly active in animated plots which use fading or transitions over time, but even in static plots, continuity can be activated when multiple time series plots are shown together. Figure 12 shows an example. Four examples are displayed as overlaid time series (top), and another four are shown as scatterplots (bottom). In the time series, plot 4 is likely perceived as having stronger association than plot 3, due to the lines moving roughly in a similar pattern. Strong negative association is not easily detected from overlaid line plots, but is easily seen in a scatterplot (Tomasetti 2015). If negative association is suspected, either using a scatterplot, or inverting one series are suggested.

Overlaying a few time series on a common scale is one way to activate the heuristic of common fate, but this does not scale well to larger numbers of simultaneous measurements. Javed et al. (2010) provide experimental evidence that small multiples are better than other alternatives when there are many simultaneous time series to display and the series cover a large visual span.

3.1.5. Working Memory. Another cognitive limitation that affects plot comprehension is the limit on working memory. Typically, working memory is limited to approximately 7 (plus or minus 2) items, or ‘chunks’. In practice, this means that categorical scales with more than 7 categories decrease readability, increase comprehension time, and require significant

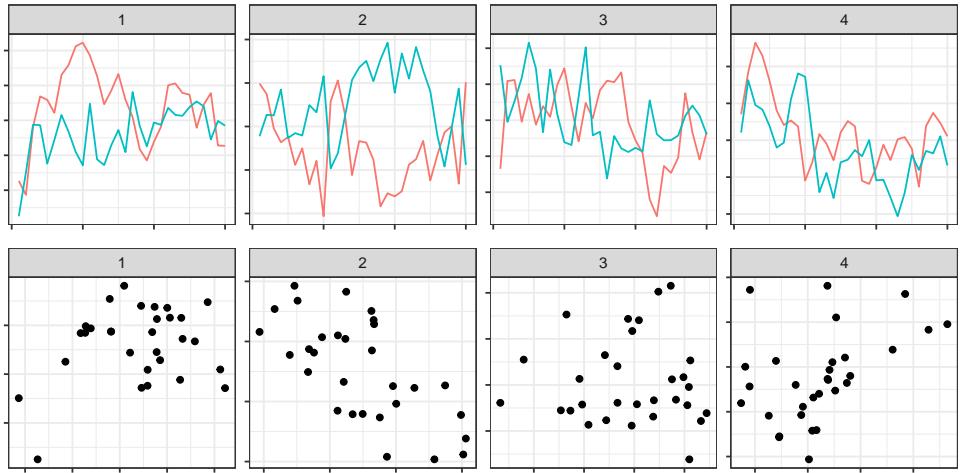


Figure 12: Four examples of pairs of time series displayed as line plots (top) and scatterplots (bottom). Perception of correlation in overlaid line plots is an example where common fate may dominate. In these four examples of line plots, viewers will likely say that plot 4 has the most similar series, but it is actually plot 2, with strong negative correlation between the two series. A scatterplot clearly shows the negative association.

attentional resources, because it is not possible to hold the legend mapping in working memory.

3.1.6. Change blindness. Not all information is stored in memory (working or long-term) as it is represented. Simons & Levin (1997) suggests that the vivid details available to us instantaneously when examining the world fade quickly, and are replaced with broad semantic descriptions or less meaningful cognitive representations of a scene. As a result of this compression, it is difficult to identify changes between similar scenes. This phenomenon, known as change blindness, affects both static and interactive plots. In static plots, it can be difficult to compare between different small multiples or facets, because the content of the plots are not reliably represented in working memory when switching attention between them. In animated plots, it is important to use transition effects to connect successive frames of the animation: this reduces change blindness and also activates the gestalt principle of common fate, allowing us to quickly identify groups of objects which are transitioning in the same direction.

3.1.7. Ease of Comparisons. Much of the psychophysics research on statistical charts examines the accuracy of comparisons and quantitative evaluations made during the process of understanding a plot. This research can be distilled into a hierarchy of comparisons (based primarily on (Cleveland & McGill 1984, 1985, Shah & Miyake 2005, Lewandowsky & Spence 1989), ranking tasks by their accuracy and difficulty as follows (roughly equivalent tasks are listed together):

1. Position (common scale)
2. Position (non-aligned scale)

3. Length, Direction, Angle, Slope
4. Area
5. Volume, Density, Curvature
6. Shading, Colour Saturation, Colour Hue
7. Discriminable Shape
8. Indiscriminable Shape

Examples of these charts are shown in Figure 13.

This ranking of cognitive tasks provides some consistent guidance for chart design: if the same data can be represented in a way that allows the user to make a comparison more accurately (based on the hierarchy), then that design is preferable. Thus, this hierarchy indicates that in most cases, it is better to use a stacked bar chart than the equivalent polar-coordinate pie chart, because a stacked bar chart requires evaluation of length, while a pie chart requires area comparisons. If information can be shown on an x or y axis rather than using color (saturation, hue, or shading), it will be easier to make numerical comparisons - we can generally order information based on color, but estimation of numerical quantities is much less precise using color than position. While the hierarchy of graphical comparisons provides some guidance, there are other design choices that can be informed by experimental research in a less systematic way.

3.2. Chart Design

VanderPlas & Hofmann (2017) provides compelling evidence that the aesthetics used in a chart can significantly affect how plots of the same data are read, and explained these differences relative to the gestalt heuristics activated by each combination of aesthetics. The use of redundant aesthetics which activate the same gestalt principles (such as colour and shape in a scatterplot, which both activate similarity) results in higher identification of corresponding data features. In addition, dual-encoding increases the accessibility of a chart to individuals who have impaired colour vision or perceptual processing (e.g. dyslexia, dysgraphia). This experimental evidence directly contradicts the guidelines popularized by Tufte 1991, which suggest the elimination of any feature which is not dedicated to representing the core data, including redundant encoding and other unnecessary graphical elements.

3.2.1. Colour. While historically there were constraints on the use of colour in graphics due to technological limitations and the economics of printing, these restrictions have evaporated with the advent of computer-generated graphics, relatively inexpensive colour printing, and an increasing tendency to share charts and graphs digitally instead of in print. Colour can encode categorical and continuous variables, and when used effectively, provides a nearly effortless way to group plot elements using a medium that does not require conscious attention. Unfortunately, while there are many ways to use colour correctly in a plot, there are generally even more ways to use it incorrectly.

Colour scales should be chosen to best match the data values and the plot type: if the goal is to show magnitude, a univariate colour scheme is typically preferable, while a double-ended colour scale is typically more effective when showing data which differs in sign and magnitude. Where possible, colour scales should use a minimal number of hues, varying intensity or lightness of the colour to show magnitude, and transitioning through neutral colours (white, light yellow) when utilizing a gradient. Figure 14 shows an exam-

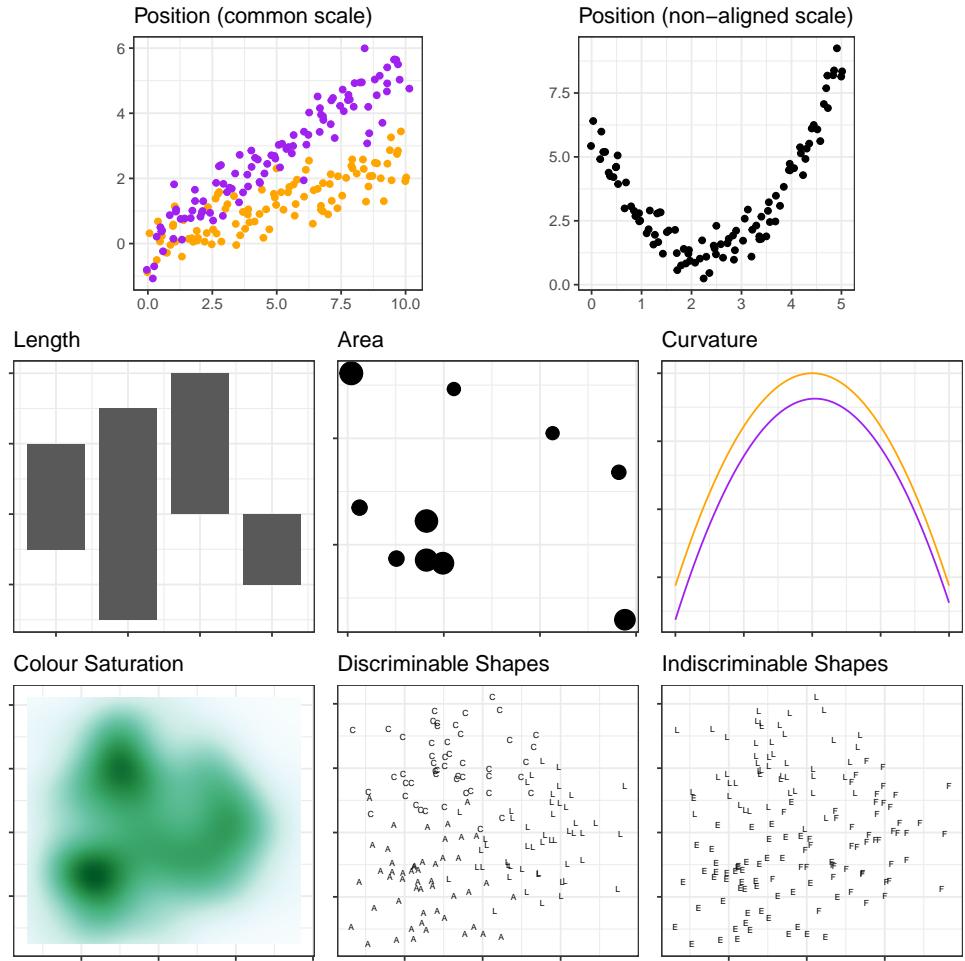


Figure 13: Examples of comparisons, by difficulty. The first plot shows two sets of points plotted on a common scale; it is easy to compare the relative trajectories of the two lines. The second plot shows another set of points which could be compared to the points in the first plot, however, because the scales are not common between the two, accurate comparisons are more difficult. In the third plot, the rectangles must be compared by length, because they are do not start or end at the same point. In the fourth plot, points are sized based on another variable, requiring comparisons of area; these comparisons are harder to make with numerical accuracy. The fifth plot shows two curves which are not completely parallel, requiring us to make a comparison of curvature; in addition, note that the difference between the two curves is hard to perceive accurately (VanderPlas & Hofmann 2015). The sixth plot shows two-dimensional density using the fill of the tiles; it is possible to make ordered comparisons here but much more difficult to estimate numerical values from the plot even when a legend is provided. The seventh and eighth plots show the same data as the sixth plot, but using discriminable shapes and indiscriminable shapes respectively. Discriminable shapes have different features, while indiscriminable shapes tend to have the same features and require more cognitive effort to separate the clusters.

Divergent Colour Schemes

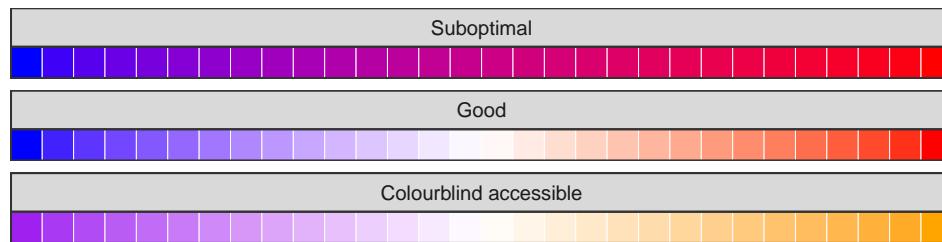


Figure 14: Examples of diverging colour schemes which are perceptually suboptimal (no transition through a neutral colour), good, and colourblind accessible. The purple - orange gradient through white is distinguishable by individuals with any type of colour deficiency affecting one of the three cones in the retina.

ple of perceptually suboptimal, good, and colourblind accessible diverging colour schemes. Cognitive load can also be reduced by selecting colours which have cultural associations that match the data display, such as the use of blue for men and red (or pink) for women, or the use of blue for cold temperatures and red/orange for warm temperatures.

It is also important to consider the human perceptual system, which does not perceive hues uniformly: we can distinguish more shades of green than any other hue, and fewer shades of yellow, so green univariate colour schemes will provide finer discriminability than other colours because the human perceptual system evolved to work in the natural world where shades of green are plentiful. Figure 15 shows the CIE 1931 colour space which maps the wavelength of a colour to a physiologically based perceptual space; a significant portion of the colour space is dedicated to greens and blues, while much smaller regions are dedicated to violet, red, orange, and yellow colours. This unevenness in mapping colour is one reason that the multi-hued rainbow colour scheme is suboptimal – the distance between points in a given colour space may not be the same as the distance between points in perceptual space (Light & Bartlein 2004, Wakita & Shimamura 2005, Borland & Ii 2007). As a result of the uneven mapping between colour space and perceptual space, multi-hued colour schemes are not recommended.

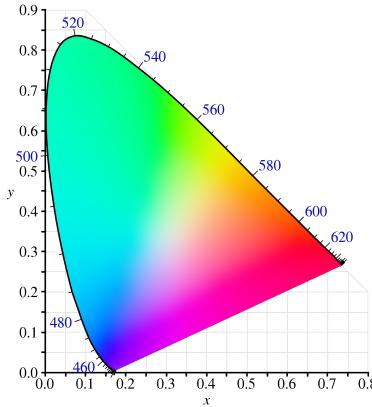


Figure 15: The CIE 1931 colour space chromaticity diagram. The CIE 1931 colour spaces defined the colour space based on physiologically perceived colours in human colour vision. The outer boundary of the curve marks the colour spectrum, with wavelengths in nm. Note that the distance between successive markings on this boundary is not constant, indicating that the perceptual distance between colours does not match the physical distance in wavelength. More broadly, there is much less area devoted to red, orange, yellow, and violet compared to blue and green; as a result, a rainbow colour scheme with equal spacing by wavelength would perceptually over-emphasize the tails of the range of values represented by the scheme.

While colour is an extremely useful aesthetic for the majority of the population, between 4 and 8 percent of males (and a much smaller proportion of females) have some sort of colour perception deficiency ('colourblindness') (Wakita & Shimamura 2005) which reduces the space of distinguishable colours. Colourblindness is common enough that it is reasonable to expect that any given chart used in a presentation or publication will be read by someone who is colour deficient. The use of dual encoding allows colourblind individuals to more readily read graphics which utilize colour, and as hue and lightness can be varied separately, it is possible to use dual-encoding without adding another aesthetic. If accessibility is the goal, default colors used in ggplot2 (Wickham 2016) should be avoided, as the saturation is held constant in this color scheme and only the hue of the colors is varied. A perceptually more successful approach to color manipulation can be found in the R package colorspace (Zeileis et al. 2009). This package is also based on the HCL (Hue, Chroma, Luminance) space of colors, but makes use of both luminance and hue when setting up color schemes and allows an assessment of colors on these three dimensions. There are several pre-existing colour schemes which may be more discriminable to colour deficient individuals, such as those provided in Lumley (2013) or Brewer (2019). Similarly, when color accessibility of a plot is a primary consideration the use of the default grey background in ggplot2 and other plotting systems should be replaced with white to maximize the contrast between background and plotted features. The default grey background for plots, which by now has become ggplot2's signature look, dates back to a recommendation in Carr (1994) to make gridlines in plots readable without dominating the data. Obviously, light grey gridlines on a white background serve the same purpose.

4. OPEN QUESTIONS AND FUTURE RESEARCH

One of the areas in graphics showing the most recent growth and creating the most excitement in the community are interactive graphics. There are a huge abundance of applications and graphics claiming to be *interactive* - yet then (Swayne & Klinke 1999) and now there is not much agreement over what 'interactive graphics' actually means. While any communication between user and device can be considered an interaction, interactive graphics should be defined as a user-driven direct manipulation of plots and plotting elements with an immediate reaction (Unwin 1999, Becker et al. 1987, Eick & Wills 1995). One of the

foundations of implementing interactive graphics is the formalization of the data pipeline necessary to support user interactions with the plot (Buja et al. 1988, Wickham et al. 2009, Lawrence et al. 2009, Xie et al. 2014). The R package `shiny` (Chang et al. 2019) is built on this idea; thousands of interactive dashboards and web-applications have been created using that platform.

Formal testing for interactive graphics is difficult without an established and generally accepted grammar of interactive graphics. However, there are several promising approaches into this direction, such as the Python framework `vega-lite` (Satyanarayan et al. 2017), and R packages, e.g. `ggvis` (Chang & Wickham 2016) and `plotly` (Sievert 2018), and most recently `ggvega` (Yang et al. 2019). Least well defined in these grammars is usually the aspect of linked brushing and highlighting in plots (Becker & Cleveland 1987), a technique crucial to interactive graphics for defining and exploring subsets of the data. An additional problem for interactive graphics is reproducibility of a user's work. By definition, interactive graphics enable a flow, rather than a single static result. Very recent advances such as `trackr` (Becker et al. 2019) not only record the user's interactions with the data, but also try to infer the user's intent by collecting metadata and automatically analysing the structure of the data and the code. This approach might be used in testing, both to evaluate different user's approaches as well as the tools used.

Another open question is the acceptance of results – a lot of historical research on best practices exists, but how much of it is being put into practice? One does not have to look far to find astonishingly bad graphics in astonishingly good outlets. We are not going to point fingers, here, but will defer to online communities such as <https://www.reddit.com/r/dataisugly/> that have passionate discussions on the worst mis-uses of graphics. What can we do about this? It is up to statisticians to teach more graphics and better graphics in the classroom in the hope to slowly change the climate. One incentive in the adoption of best practices might be where there is a cost benefit, such as when assessing business performance, or a marketing campaign, there is scope for providing a measure on which to gauge effective visual communication. There is also significant room to improve plots in academic publications in order to better communicate research results (Unwin in press).

In the 100 years of empirical evaluation of the perception and utility of statistical graphics, we have assembled a working knowledge of how to best create graphs which are easily read and understood. Once rare, charts are now everywhere we look - on the news, in papers and magazines, and online in interactive form. Going forward, we must do a better job of translating the academic research into practice, making it easier for academics and non-academics alike to create useful, well-designed graphics.

DISCLOSURE STATEMENT

The authors are not aware of any affiliations, memberships, funding, or financial holdings that might be perceived as affecting the objectivity of this review.

ACKNOWLEDGMENTS

Heike Hofmann and Susan Vanderplas are partially supported by the Center for Statistics and Applications in Forensic Evidence (CSAFE) through Cooperative Agreement #70NANB15H176 between NIST and Iowa State University, which includes activities carried out at Carnegie Mellon University, University of California Irvine, and University of

Virginia.

LITERATURE CITED

- Amer T. 2005. Bias due to visual illusion in the graphical presentation of accounting information. *Journal of Information Systems* 19:1–18
- Baird JC, Lewis C, Romer D. 1970. Relative frequencies of numerical responses in ratio estimation. *Perception & Psychophysics* 8:358–362
- Becker G, Moore SE, Lawrence M. 2019. trackr: A Framework for Enhancing Discoverability and Reproducibility of Data Visualizations and Other Artifacts in R. *Journal of Computational and Graphical Statistics* :1–15
- Becker RA, Cleveland WS. 1987. Brushing Scatterplots. *Technometrics* 29:127
- Becker RA, Cleveland WS, Wilks AR. 1987. Dynamic graphics for data analysis. *Statist. Sci.* 2:355–383
- Beecham R, Dykes J, Meulemans W, Slingsby A, Turkay C, Wood J. 2017. Map LineUps: Effects of spatial structure on graphical inference. *IEEE Transactions on Visualization and Computer Graphics* 23:391–400
- Beniger JR, Robyn DL. 1978. Quantitative Graphics in Statistics: A Brief History. *The American Statistician* 32:1–11
- Bernstein EM, Cowden DJ. 1937. Graphic Presentation of Trend Data. *Southern Economic Journal* 3:443–451
- Bertin J, Berg WJ. 1983. Semiology of graphics: diagrams, networks, maps, vol. 1. University of Wisconsin press Madison
- Borland D, Ii RMT. 2007. Rainbow Color Map (Still) Considered Harmful. *IEEE Computer Graphics and Applications* 27:14–17
- Brandes D. 1976. The Present State of Perceptual Research in Cartography. *The Cartographic Journal* 13:172–176
- Brewer CA. 2019. Color Ddvice for Cartography. <http://colorbrewer2.org>
- Brinton WC. 1917. Graphic methods for presenting facts. Engineering magazine company
- Broersma H, Molenaar IW. 1985. Graphical perception of distributional aspects of data. *Computational Statistics Quarterly* 2:53–72
- Buja A, Asimov D, Hurley C, McDonald JA. 1988. Elements of a viewing pipeline for data analysis. In *Dynamic Graphics for Statistics*, eds. WS Cleveland, ME McGill, chap. 11. Wadsworth Inc, 277–308
- Buja A, Cook D, Hofmann H, Lawrence M, Lee EK, et al. 2009. Statistical inference for exploratory data analysis and model diagnostics. *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences* 367:4361–4383
- Callaghan TC. 1984. Dimensional interaction of hue and brightness in preattentive field segregation. *Perception & Psychophysics* 36:25–34
- Carr DB. 1994. Using gray in plots. *Statistics Computing and Graphics Newsletter* 5
- Carswell CM, Wickens CD. 1987. Information integration and the object display An interaction of task demands and display superiority. *Ergonomics* 30:511–527
- Chandar N, Collier D, Miranti P. 2012. Graph standardization and management accounting at AT&T during the 1920s. *Accounting History* 17:35–62
- Chang W, Cheng J, Allaire J, Xie Y, McPherson J. 2019. shiny: Web application framework for r. R package version 1.3.2
- Chang W, Wickham H. 2016. ggvis: Interactive grammar of graphics. R package version 0.4.3
- Cleveland WS, McGill R. 1984. Graphical perception: Theory, experimentation, and application to the development of graphical methods. *Journal of the American statistical association* 79:531–554
- Cleveland WS, McGill R. 1985. Graphical Perception and Graphical Methods for Analyzing Scientific Data. *Science* 229:828–833
- Cleveland WS, McGill R. 1987. Graphical Perception: The Visual Decoding of Quantitative Information on Graphical Displays of Data. *Journal of the Royal Statistical Society. Series A (General)* 150:192

- Cooke L. 2010. Assessing Concurrent Think-Aloud Protocol as a Usability Test Method: A Technical Communication Approach. *IEEE Transactions on Professional Communication* 53:202–215
- Croxton FE. 1932. Graphic Comparisons by Bars, Squares, Circles, and Cubes. *Journal of the American Statistical Association* 27:54–60
- Croxton FE, Stryker RE. 1927. Bar Charts Versus Circle Diagrams. *Journal of the American Statistical Association* 22:473–482
- Desnoyers L. 2011. Toward a Taxonomy of Visuals in Science Communication. *Technical Communication* 58:16
- Dunbar K. 1995. How scientists really reason: Scientific reasoning in real-world laboratories. *The nature of insight* 18:365–395
- Dunn R. 1988. Framed Rectangle Charts or Statistical Maps with Shading: An Experiment in Graphical Perception. *The American Statistician* 42:123
- Eells WC. 1926. The Relative Merits of Circles and Bars for Representing Component Parts. *Journal of the American Statistical Association* 21:119–132
- Eick SG, Wills GJ. 1995. High interaction graphics. *European Journal of Operational Research* 81:445–459
- Fabrikant SI, Hespanha SR, Hegarty M. 2010. Cognitively Inspired and Perceptually Salient Graphic Displays for Efficient Spatial Inference Making. *Annals of the Association of American Geographers* 100:13–29
- Fienberg SE. 1979. Graphical Methods in Statistics. *The American Statistician* 33:165–178
- Flury B, Riedwyl H. 1988. Multivariate statistics: A practical approach. London: Chapman & Hall
- Funkhouser HG. 1937. Historical Development of the Graphical Representation of Statistical Data. *Osiris* 3:269–404
- Gelman A, Unwin A. 2013. Infovis and Statistical Graphics: Different Goals, Different Looks. *Journal of Computational and Graphical Statistics* 22:2–28
- Goldberg JH, Helfman JI. 2010. Comparing information graphics: a critical look at eye tracking, In *Proceedings of the 3rd BELIV'10 Workshop on BEyond time and errors: novel evaLuation methods for Information Visualization - BELIV '10*. Atlanta, Georgia: ACM Press
- Green TM, Fisher B. 2011. The Personal Equation of Complex Individual Cognition during Visual Interface Interaction. In *Human Aspects of Visualization*, eds. A Ebert, A Dix, ND Gershon, M Pohl, vol. 6431. Berlin, Heidelberg: Springer Berlin Heidelberg, 38–57
- Guan Z, Lee S, Cuddihy E, Ramey J. 2006. The validity of the stimulated retrospective think-aloud method as measured by eye tracking, In *Proceedings of the SIGCHI conference on Human Factors in computing systems - CHI '06*. ACM Press
- Harms H. 1991. August Friedrich Wilhelm Crome (1753-1833). *Cartographica Helvetica* 33:33–38
- Hasanhodzic J, Lo AW, Viola E. 2010. Is It Real, or Is It Randomized?: A Financial Turing Test. *arXiv:1002.4592 [cs, q-fin]* 00008 arXiv: 1002.4592
- Healey CG, Booth KS, Enns JT. 1996. High-speed visual estimation using preattentive processing. *ACM Transactions on Computer-Human Interaction (TOCHI)* 3:107–135
- Healey CG, Enns JT. 1999. Large datasets at a glance: Combining textures and colors in scientific visualization. *IEEE transactions on visualization and computer graphics* 5:145–167
- Heer J, Bostock M. 2010. Crowdsourcing graphical perception: using mechanical turk to assess visualization design, In *Proceedings of the SIGCHI conference on human factors in computing systems*. ACM
- Hofmann H, Follett L, Majumder M, Cook D. 2012. Graphical tests for power comparison of competing designs. *IEEE Transactions on Visualization and Computer Graphics* 18:2441–2448
- Hughes BM. 2001. Just Noticeable Differences in 2d and 3d Bar Charts: A Psychophysical Analysis of Chart Readability. *Perceptual and motor skills* 92:495–503
- International Statistical Congress. 1858. Emploi de la cartographie et de la méthode graphique en général pour les besoins spéciaux de la statistique. Tech. rep., Compte rendu de la troisième session du Congrès International de Statistique réuni à Vienne les 31 Août, 1, 2, 3, 4, 5, Vienne

- Javed W, McDonnel B, Elmquist N. 2010. Graphical Perception of Multiple Time Series. *IEEE Transactions on Visualization and Computer Graphics* 16:927–934
- Karsten K. 1923. Charts and Graphs: An Introduction to Graphic Methods in the Control and Analysis of Statistics. Prentice-Hall, Incorporated
- Kibirige H. 2017. plotnine: a grammar of graphics for python
- Kim S, Lombardino LJ, Cowles W, Altmann LJ. 2014. Investigating graph comprehension in students with dyslexia: An eye tracking study. *Research in Developmental Disabilities* 35:1609–1622
- Kirschenbaum SS. 2003. Comparative Cognitive Task Analysis: The Cognition of Weather Forecasting, In *Proceedings of the Human Factors and Ergonomics Society Annual Meeting*, vol. 47
- Kostelnick C. 2016. The Re-Emergence of Emotional Appeals in Interactive Data Visualization. *Technical Communication* 63:116–135
- Kruskal W. 1977. Visions of maps and graphs, In *In Proceedings of the International Symposium on Computer-Assisted Cartography, Auto-Carto II*
- Lawrence M, Wickham H, Cook D, Hofmann H, Swayne DF. 2009. Extending the ggobi pipeline from r: Rapid prototyping of interactive visualizations. *Computational Statistics* 24:195–205. Special issue for Proceedings of the 5th International Workshop on Directions in Statistical Computing.
- Legge GE, Gu Y, Luebker A. 1989. Efficiency of graphical perception. *Perception & Psychophysics* 46:365–374
- Lewandowsky S, Spence I. 1989. The perception of statistical graphs. *Sociological Methods & Research* 18:200–242
- Light A, Bartlein PJ. 2004. The end of the rainbow? Color schemes for improved data graphics. *Eos, Transactions American Geophysical Union* 85:385–391
- Loy A, Follett L, Hofmann H. 2016. Variations of Q-Q Plots: The Power of Our Eyes! *The American Statistician* 70:202–214
- Loy A, Hofmann H. 2015. Are You Normal? The Problem of Confounded Residual Structures in Hierarchical Linear Models. *Journal of Computational and Graphical Statistics* 24:1191–1209
- Lumley T. 2013. dichromat: Color schemes for dichromats. R package version 2.0-0
- MacDonald-Ross M. 1977. How Numbers Are Shown: A Review of Research on the Presentation of Quantitative Data in Texts. *AV Communication Review* 25:359–409
- Majumder M, Hofmann H, Cook D. 2013. Validation of visual statistical inference, applied to linear models. *Journal of the American Statistical Association* 108:942–956
- McDonald L. 2014. Florence Nightingale, statistics and the Crimean War. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 177:569–586
- Morel P. 2018. Gramm: grammar of graphics plotting in Matlab. *The Journal of Open Source Software* 3:568
- Netzel R, Vuong J, Engelke U, O'Donoghue S, Weiskopf D, Heinrich J. 2017. Comparative eye-tracking evaluation of scatterplots and parallel coordinates. *Visual Informatics* 1:118–131
- Normand MP, Bailey JS. 2006. The Effects of Celeration Lines on Visual Data Analysis. *Behavior Modification* 30:295–314
- North C. 2006. Toward measuring visualization insight. *IEEE Computer Graphics and Applications* 26:6–9
- Okan Y, Galesic M, Garcia-Retamero R. 2016. How People with Low and High Graph Literacy Process Health Graphs: Evidence from Eye-tracking. *Journal of Behavioral Decision Making* 29:271–294
- Peterson LV, Schramm W. 1954. How Accurately Are Different Kinds of Graphs Read? *Audio Visual Communication Review* 2:178–189
- Peuchet J, Gilbert C. 1805. Statistique élémentaire de la France. Chez Gilbert et compagnie
- Playfair W. 1801. The statistical breviary; shewing the resources of every state and kingdom in Europe. J. Wallis
- Roy Chowdhury N, Cook D, Hofmann H, Majumder M, Lee EK, Toth AL. 2015. Using visual statistical inference to better understand random class separations in high dimension, low sample

- size data. *Computational Statistics* 30:293–316
- Ryu YS, Yost B, Convertino G, Chen J, North C. 2003. Exploring Cognitive Strategies for Integrating Multiple-View Visualizations. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 47:591–595
- Satyanarayan A, Moritz D, Wongsuphasawat K, Heer J. 2017. Vega-Lite: A Grammar of Interactive Graphics. *IEEE Transactions on Visualization and Computer Graphics* 23:341–350
- Shah P, Carpenter PA. 1995. Conceptual limitations in comprehending line graphs. *Journal of Experimental Psychology: General* 124:43
- Shah P, Miyake A. 2005. The Cambridge handbook of visuospatial thinking. Cambridge University Press
- Sievert C. 2018. plotly for r
- Simons DJ, Levin DT. 1997. Change blindness. *Trends in Cognitive Sciences* 1:261–267
- Smith CD. 1996. Imago Mundi's Logo the Babylonian Map of the World. *Imago Mundi* 48:209–211
- Spence I. 1990. Visual psychophysics of simple graphical elements. *Journal of Experimental Psychology: Human Perception and Performance* 16:683
- Swayne D, Klinke S. 1999. Introduction to the special issue on interactive graphical data analysis: what is interaction? *Computational Statistics* 14:1–6
- Tan JK. 1994. Human processing of two-dimensional graphics: Information-volume concepts and effects in graph-task fit anchoring frameworks. *International Journal of Human-Computer Interaction* 6:414–456
- Teghtsoonian M. 1965. The judgment of size. *The American journal of psychology* 78:392–402
- The British Museum. 1882. The map of the world. https://www.britishmuseum.org/research/collection_online/collection_object_details.aspx?assetId=404485001&objectId=362000&partId=1
- Tomasetti N. 2015. Comparing the power of plot designs to reveal correlation. https://github.com/dicook/lineplots_v_scatterplot
- Trafton GJ, Kirschenbaum SS, Tsui TL, Miyamoto RT, Ballas JA, Raymond PD. 2000. Turning pictures into numbers: extracting and generating information from complex visualizations. *International Journal of Human-Computer Studies* 53:827–850
- Treisman AM. 1980. A Feature-Integration Theory of Attention. *Cognitive Psychology* 12:97–136
- Trickett SB, Fu WT, Schunn CD, Trafton JG. 2000a. From Dipsy-Doodles to Streaming Motions: Changes in Representation in the Analysis of Visual Scientific Data, In *Proceedings of the Annual Meeting of the Cognitive Science Society*, vol. 22
- Trickett SB, Trafton JG, Schunn CD. 2000b. Blobs, dipsy-doodles and other funky things: Framework anomalies in exploratory data analysis, In *Proceedings of the Annual Meeting of the Cognitive Science Society*, vol. 22
- Tufte E. 1991. The Visual Display of Quantitative Information. USA: Graphics Press, 2nd ed.
- Tukey J. 1972. Some Graphic and Semigraphic Displays. In *Statistical Papers in Honor of George W. Snedecor*, ed. T Bancroft. Ames, IA: The Iowa State University Press, 293–316
- Tukey JW. 1965. The Technical Tools of Statistics. *The American Statistician* 19:23–28
- Ulery BT, Hicklin RA, Buscaglia J, Roberts MA. 2011. Accuracy and reliability of forensic latent fingerprint decisions. *Proceedings of the National Academy of Sciences* 108:7733–7738
- Unwin A. 1999. Requirements for interactive graphics software for exploratory data analysis. *Computational Statistics* 14:7–22
- Unwin A. in press. Why is data visualization important? what is important in data visualisation? *The Harvard Data Science Review*
- VanderPlas S, Goluch R, Hofmann H. 2019. Framed! Reproducing and Revisiting 150 year old charts. *Journal of Computational and Graphical Statistics*
- VanderPlas S, Hofmann H. 2015. Signs of the Sine Illusion—Why We Need to Care. *Journal of Computational and Graphical Statistics* 24:1170–1190
- VanderPlas S, Hofmann H. 2017. Clusters Beat Trend!? Testing Feature Hierarchy in Statistical

- Graphics. *Journal of Computational and Graphical Statistics* 26:231–242
- von Huhn R. 1927. Further Studies in the Graphic Use of Circles and Bars: A Discussion of the Ell's Experiment. *Journal of the American Statistical Association* 22:31–39
- Wagemans J, Elder JH, Kubovy M, Palmer SE, Peterson MA, et al. 2012a. A Century of Gestalt Psychology in Visual Perception: I. Perceptual Grouping and Figure-Ground Organization. *Psychological Bulletin* 138:1172–1217
- Wagemans J, Feldman J, Gepshtein S, Kimchi R, Pomerantz JR, et al. 2012b. A Century of Gestalt psychology in Visual Perception: II. Conceptual and Theoretical Foundations. *Psychological Bulletin* 138:1218–1252
- Wakita K, Shimamura K. 2005. SmartColor: disambiguation framework for the colorblind, In *Proceedings of the 7th international ACM SIGACCESS conference on Computers and accessibility - Assets '05*. Baltimore, MD, USA: ACM Press
- Walker FA. 1874. Statistical Atlas of the United States, Based on the Results of the Ninth Census, 1870, with Contributions from Many eminent Men of Science, and Several Departments of the Government. [New York] J. Bien, lith
- Wickham H. 2010. A Layered Grammar of Graphics. *Journal of Computational and Graphical Statistics* 19:3–28
- Wickham H. 2013. Graphical criticism: some historical notes. *Journal of Computational and Graphical Statistics* 22:38–44
- Wickham H. 2016. ggplot2: Elegant graphics for data analysis. Springer-Verlag New York
- Wickham H, Chowdhury NR, Cook D, Hofmann H. 2018. nullabor: Tools for graphical inference. R package version 0.3.5
- Wickham H, Cook D, Hofmann H, Buja A. 2010. Graphical inference for infovis. *IEEE Transactions on Visualization and Computer Graphics* 16:973–979
- Wickham H, Grolemund G. 2017. R for data science: Import, tidy, transform, visualize, and model data. O'Reilly Media, Inc., 1st ed.
- Wickham H, Lawrence M, Cook D, Buja A, Hofmann H, Swayne DF. 2009. The plumbing of interactive graphics. *Computational Statistics* 24:207–215. Special issue for Proceedings of the 5th International Workshop on Directions in Statistical Computing
- Widen HM, Elsner JB, Pau S, Uejio CK. 2016. Graphical inference in geographical research. *Geographical Analysis* 48:115–131
- Wilkinson L. 1999. The Grammar of Graphics. Statistics and computing. Springer
- Woller-Carter MM, Okan Y, Cokely ET, Garcia-Retamero R. 2012. Communicating and Distorting Risks with Graphs: An Eye-Tracking Study. *Proceedings of the Human Factors and Ergonomics Society Annual Meeting* 56:1723–1727
- Xie Y, Hofmann H, Cheng X. 2014. Reactive programming for interactive graphics. *Statist. Sci.* 29:201–213
- Yang W, Jeppson H, Lyttle IJ. 2019. ggvega: Translator from 'ggplot2' to 'vega-lite'. R package version 0.0.0.9001
- Yates J. 1985. Graphs as a Managerial Tool: A Case Study of Du Pont's Use of Graphs in the Early Twentieth Century. *Journal of Business Communication* 22:5–33
- Zeileis A, Hornik K, Murrell P. 2009. Escaping RGBland: Selecting colors for statistical graphics. *Computational Statistics & Data Analysis* 53:3259–3270
- Zgraggen E, Zhao Z, Zeleznik R, Kraska T. 2018. Investigating the Effect of the Multiple Comparisons Problem in Visual Analysis, In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems - CHI '18*. Montreal QC, Canada: ACM Press