

DATA SCIENCE, DASHBOARDS, AND THE WAY IT WORKS WITH
STATISTICS

by

Denise Renee Bradford

A DISSERTATION

Presented to the Faculty of
The Graduate College at the University of Nebraska
In Partial Fulfillment of Requirements
For the Degree of Doctor of Philosophy

Major: Statistics

Under the Supervision of Susan R. VanderPlas, Ph.D

Lincoln, Nebraska

Month, Year

DATA SCIENCE, DASHBOARDS, AND THE WAY IT WORKS WITH
STATISTICS

Denise Renee Bradford, Ph.D.

University of Nebraska, Year

Adviser: Susan R. VanderPlas, Ph.D

Here is my abstract. (*350 word limit*)

COPYRIGHT

© Year, Denise Renee Bradford

DEDICATION

Dedicated to...

ACKNOWLEDGMENTS

Thank you to all my people!

Table of Contents

List of Figures	vii
List of Tables	viii
1 Literature Review	1
1.1 Introduction	1
1.2 History of Exploratory Data Analysis (EDA)	3
1.3 History/Review of Graphical Representation	8
1.4 Dashboard Design/ Human Perception (Visual Information) .	16
1.5 Conclusion	30
2 Chapter Paper on Rural Shrink Smart Manuscript submitted to Journal of Data Science Special Issue	31
2.1 Abstract	31
2.2 Introduction	31
2.3 Data Description	33
2.4 Dashboard Design Considerations	35
2.5 Guiding Design Principles	38
2.6 Dashboard Design Process	39
2.7 Discussion	47
2.8 Future Work	48

2.9 Conclusions	50
3 Chapter 2 Stuff	51
4 Tables, Graphics, References, and Labels	52
4.1 Tables	52
4.2 Figures	54
4.3 Footnotes and Endnotes	58
4.4 Cross-referencing chapters and sections	59
4.5 Bibliographies	61
4.6 Anything else?	63
Conclusion	64
A The First Appendix	65
B The Second Appendix, for Fun	67
Colophon	68
References	73

List of Figures

1.1	Dashboard Design Mind Map	1
1.2	Roles of Graphics	2
1.3	Grammar of Graphics Diagram of Wickham and Wilkinson's work	6
1.4	Effective Data Graphics	7
1.5	Big Data Diagram	7
1.6	van Wijj Simple Visualization Model	11
1.7	10 Type Tasks with Sources	13
1.8	Interaction Types with Sources	14
1.9	Cognitive Structures Citation Timeline	19
1.10	Visual Reference model by Card	21
1.11	Elements of EDA Workflow	27
2.1	Diagram of considerations for our dashboard design process.	35
2.2	Initial dashboard design sketch (top) and implementation (bottom).	42
2.3	A second iteration of the sketched design (top) and the implementation (bottom).	45
4.1	logo	55
4.2	Mean Delays by Airline	56
4.3	Subdiv. graph	58
4.4	A Larger Figure, Flipped Upside Down	59

List of Tables

4.1 Correlation of Inheritance Factors for Parents and Child	53
------------------------------------------------------------------------	----

Chapter 1

Literature Review

1.1 Introduction

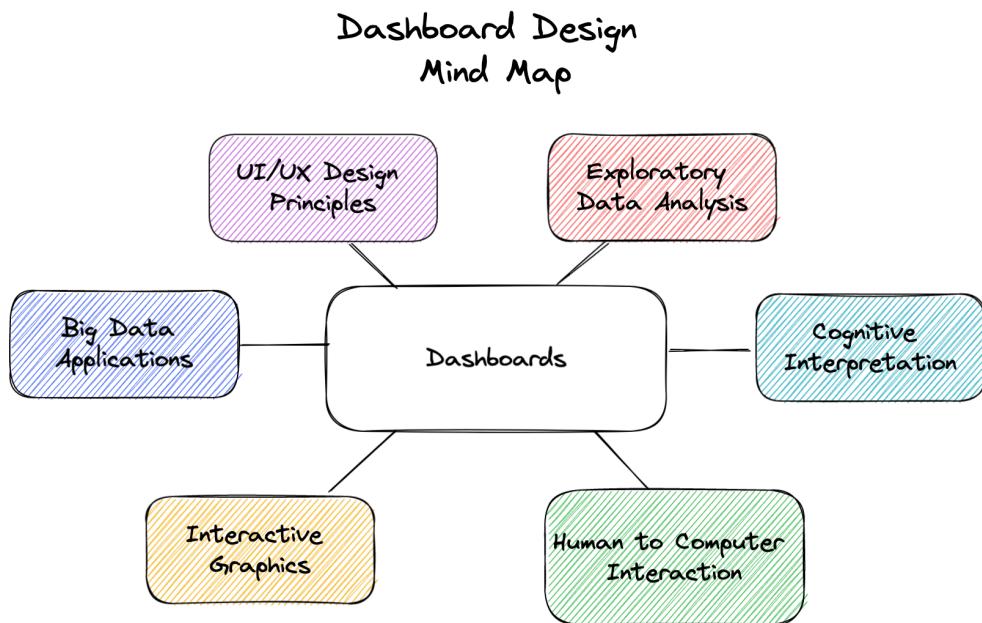


Figure 1.1: Dashboard Design Mind Map

Dashboard Design with naive users at the center of the design with adding new graphical visualization and testing the effectiveness of new graphics.

Data Visualizations are described as a graphical representation of tabular data and information. Data visualization tools are used in ways to accessible

way to see and understand trends, outlines, and patterns in data.

The literature review will explore the facets of Exploratory Data Analysis (EDA) through the history of graphical representation, followed by the study of Dashboard Design, Visual Information, and Big Data Analytics.

Roles of graphics in Data Analysis: Graphics and Tables are forms of communication that we are looking to determine:

What is the audience? Our audience is those users who “know” their data but cannot utilize it due to a lack of skills and resources.

What is the message? The message will be to understand how to communicate visual information of statistical graphics in an ethical and understandable way.

We can break down statistical graphics into two categories: - Analysis: design to see patterns, trends, and the process of data description, interpretation - Presentation: design to attract attention, make a point, illustrate a conclusion

Based on the following diagram

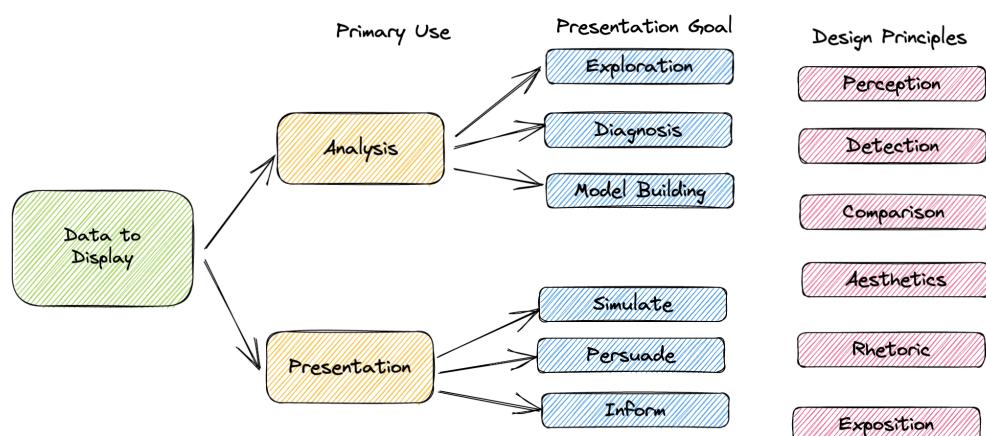


Figure 1.2: Roles of Graphics

As tabular data grows in our society, it has pushed our need to create visualizations that represent that data in a correct and digestible way. The work that has been done in statistical graphics research would suggest that the work that is done in this work has grown increasingly in formal adding sound guidelines with statistically sound metrics that are needed. This work has been explored in various disciplines but not in statistical graphics. We know that dashboards are visual information tools that include and have multiple statistical pictures. Thus, we should consider and think about the fact that we may need to consider the statistical and psychological framework in a dashboard that will combine multiple statistical graphics.

One graphic can be misleading alone, but can two graphics be considered misleading when placing two “correctly” designed graphics together on a dashboard? This is something that visual information researchers explore but don’t test. We would like to test the combination with interactive features, such as map selection and highlighted lines in Parallel Coordinate Plot, for example.

Since we don’t completely understand the way that basic statistical graphics are

When effectively adding advanced graphical representations to a naive user, what have we done to create a proper understanding of an engaging idea of what’s going on?

1.2 History of Exploratory Data Analysis (EDA)

Visualizations of data are essential for exploratory data analysis (EDA) along with model diagnostics. Plots for EDA are a valuable tool for guiding an analyst in discovering the relationships between variables in their data.

When using plots in model diagnostics, plots help analysts determine whether or not the model is an appropriate way to model. During the initial EDA stage, an analyst may find that a variable or a covariate is directly related to the dependent variable when looking at a correlation heatmap or a scatterplot. This will be important to know before starting a linear model analysis. Much of our general understanding is from introductory statistics courses. The basic understanding can be formalized to visualize the discovery process.

1.2.1 Study 1: EDA through the eyes of Tukey

John Tukey was the first to organize the collection and methods associated with philosophy into Exploratory Data Analysis (EDA). Tukey was a person who helped utilize the likes of stem-and-leaf plot, boxplot-resistant smooth, and rootgram.

Tukey's Principles in EDA:

1. Graphical exploration looking for patterns or displaying fit.
 - The method demonstrates things about data that are not understood by a single numeric metric. This has been useful in graphing the data before you develop summary statistics.
2. Describing the general patterns of the data.
 - This step should be insensitive to outliers. In general, think about the types of resistant measures (i.e., median or mean). This step is making sure to determine data patterns.

3. The natural scale/state that the data are at their best. This will be the step at which the scale of data can be helpful for analysis. The reexpressing data to a new scale by taking the square root or logarithmic scale.
4. The mostly known parts of EDA but is done in the way of accessing fit of the data. This is taught in every statistics 101 class. The growth of machine learning and prediction methods have now used residuals more in the toolbox to access the best prediction models.
 - The idea generally is to determine the deviations in the data from a general pattern by looking at the data from the fit of the data.

1.2.2 EDA through the eyes of the modern grammar of graphics

The grammar of graphics (gg of ggplot2) is a theory that is well-defined for creating statistical graphics with work from Wilkinson (**Wilkinson1999?**) and Hadley Wickham (Wickham, 2016). Leland Wilkinson created the structure and defined the phrase. Grammar of graphics is defined as the framework which follows a layered approach to describe and construct visualizations or graphics in a structured manner.

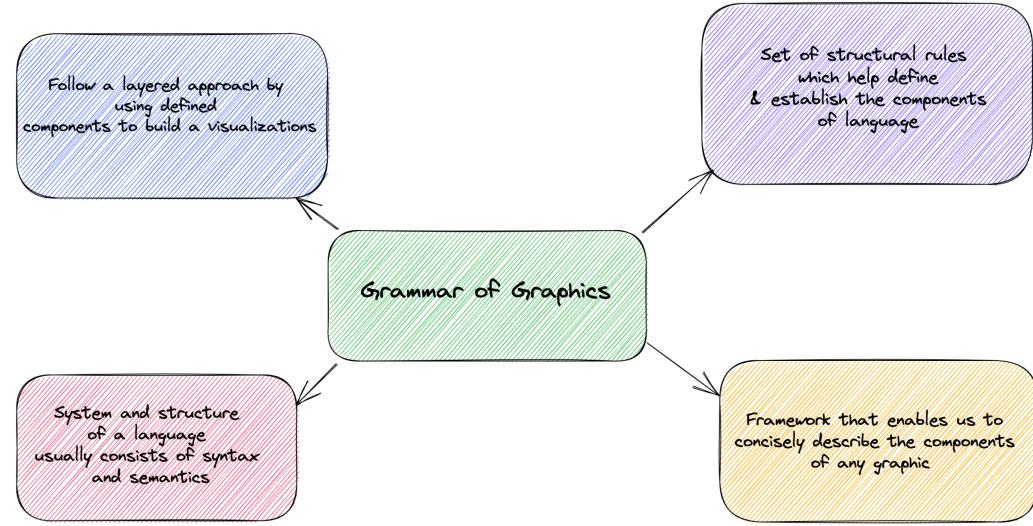


Figure 1.3: Grammar of Graphics Diagram of Wickham and Wilkinson's work

The theory of graphics is based on plot layers, which are built with four distinct pieces:

- data - *tabular formatted from user*
- aesthetic mapping - *takes data and maps the variable in the data frame to a particular visual features*
- statistical transformation - *determines how to transform the data to values to create the visual feature*
- geometric object and a position adjustment - *geom object is used to draw a plot layer, and position adjustment is to help with adjusting the visual feature in the space*

Seven major components help build effective visualizations that build on one another. Seven major components help build effective visualizations. The following diagram will display the effectiveness of data graphics:

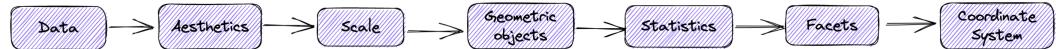


Figure 1.4: Effective Data Graphics

1.2.3 EDA through the eyes of Big Data and Big Data Analytics

Big Data Analytics is where advanced analytic techniques are applied to big data sets. - Analytics based on extensive data samples reveals and leverages business change. The larger the data set, the more difficult it becomes to manage. This paradigm has only become more and more of the pieces of data exploration since tabular data is growing.

1.2.3.1 Characteristics of Big Data

Big Data Characteristics



Figure 1.5: Big Data Diagram

The Big Data Diagram will display the three main features and characteristics of big data: - Volume: is its size and how enormous it is - Variety: includes the different formats and types of data, as well as the different kinds of uses and ways of analyzing the data - Velocity: refers to the rate with which data is changing or how often it is created

1.2.3.2 Big Data Storage and Management

We can't discuss big data analytics without discussing data storage and management. Our tasks with big data storage is to determine where and

how this data will be stored once it is acquired:

- Traditional structured data storage and retrieval methods include relational databases, data marts, and data warehouses.
- Uploaded to the storage from operational data stores using: Extract, Transform, Load (ETL) or Extract, Load, Transform (ELT) tools which extract data from outside sources, transform the data to fit operational needs, and transform and cataloged before made available for data mining and online analytical functions.

Given the growing numbers of data sources and the sophistication of data analyses, big data storage should allow analysts to produce and adapt data rapidly.

- This requires an agile database whose current data analyses use complex statistical methods, and analysts need to be able to be deep and serve as a sophisticated algorithmic run time engine.

1.3 History/Review of Graphical Representation

1.3.1 History/Origin of Statistical Visualization

Visualization is NOT prediction. Much of the field is focus on prediction models but big data analytics visualization is an important part of working to develop methodology for helping analysts to explore and understand what patterns are presented.

The data from plots can be very different from what would be learned by modeling and prediction which means both types of summarization are equally important.

Visual inference uses our ability to detect graphical anomalies. The idea of formal testing remains the same in visual inference – with one exception: The test statistic is now a graphical display which is compared to a “reference distribution” of plots showing the null.

1.3.2 Statistical Visualization Timeline

Categorical Data Visualization Friendly(Friendly, 2014) detailed using SAS with hand-on experiment to visually present categorical data analysis. PCPs have been used by researchers, to visualize categorical data. Beygelzimer, Perng, Ma (Beygelzimer, Perng, & Ma, 2001); Ma and Hellerstien (Ma & Hellerstein, 2001) created a fast ordering categorical data analysis algorithm helped visualization, where their algorithms help organize the original parallel coordinate plots clearer. Hammock plots are modified versions of parallel coordinate plots, which was invented by Schonlau to visualize categorical data (Schonlau, 2003). His design replace coordinate polygons by rectangles to present the number. Treemaps are modified to support categorical data visualization.

CatTree gives a hierarchical categorical data visualization with interaction (Kolatch & Weinstein, 2001). Fernstad developed an interactive system combining parallel coordinate, tables and scatterplot matrix together for an overview explorative analysis. Fully research on categorical data visualization to support algorithm understanding. Novel contingency wheel presented by (Alsallakh, Gröller, Miksch, & Suntiger, 2011) support visual analytics in categorical data and he measured association based on Pearson’s residuals and use visual abstraction based on elements frequency.

High-Dimensional Data Visualization A popular research area in visualization research since high-dimensional data is always fuzzy to mining. Direct visualization includes geometric visualizations:

- scatterplots: simply use dot in coordinate to present data point.
- parallel coordinates: present each dimension as axes and every data item intersects dimensions as polygon line at certain position
- RadViz/ PolyViz
- GridViz

Besides these traditional geometric visualization methods, iconographic displays like human faces and star glyphs used funny ways to present multivariate data. Hierarchical methods used widely in parallel coordinates, which give analysts an intuitive view of clustering information (Fua, Ward, & Runnensteiner, 1999); (Johansson, Ljung, Jern, & Cooper, 2005). Rearrange the dimensions by dimensions similarity on parallel coordinates, circle segments and recursive patterns. (Ankerst, Berchtold, & Keim, 1998). (Guo, 2003) used interactive feature selection method help users identify interesting subspaces from high-dimensional data sets.

Visualizing Complex Data

William Cleveland's subcycle plots: - glyphmaps and binned graphics that are emerging from big data visualization efforts.

- glyphs and other plots have been embedded in maps

Bertin's semiologic of Graphics, a seminal work in the academic study of visualization. Glyphmaps have been developed as a tool for tracking climate

and climate change data (Wickham, Hofmann, Wickham, & Cook, 2012); (**hobbs2010?**)

Interactive Graphics

The area of interactive graphics is still very much a work in progress despite existing as field of research since the late 1960s. Developments driven in part by new technology, such as d3 (Bostock, Ogievetsky, & Heer, 2011).

Visualizations are more than just a picture. They are now a tool that facilitates analytic activity through different modes of interaction. (Yi, Kang, Stasko, & Jacko, 2007). The term visualization is context-free, as it can mean different things to different people depending on the situation (Parsons & Sedig, 2014). van Wijj's (Van Wijk, 2005) simple visualization model shows how insights are generated as the human participates in a feedback loop between reading and interacting with visualization. This model is also context-free allow for the focus to be on the feedback loops between visualization and user.

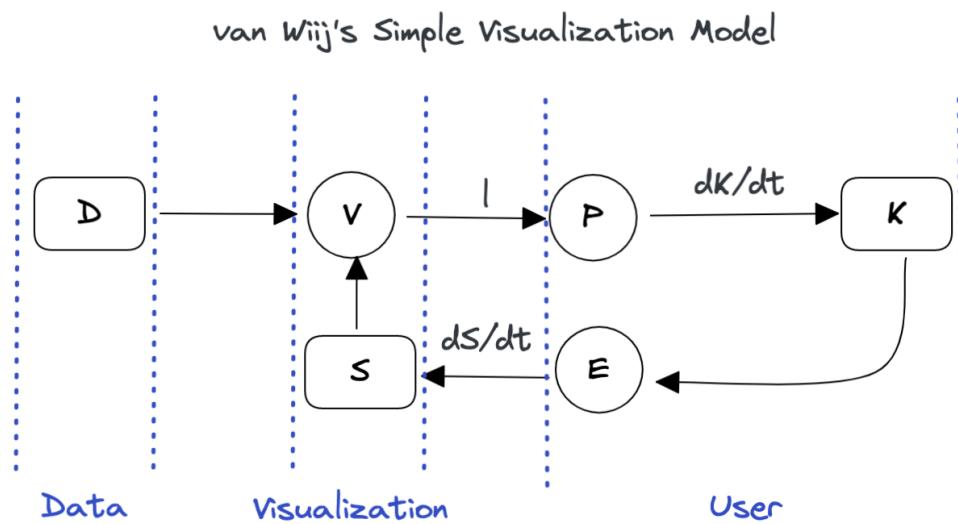


Figure 1.6: van Wijj Simple Visualization Model

Interaction allows the user to define what data they see and how they see the data creating a dialogue between the user and the system. Theories behind visual representation include: - graphical comprehension ((Cleveland & McGill, 1984)) - preattentive processing ((Ware, 2012)) - gestalt theory ((Few, 2009)) - graphical excellence ((**tufte2001?**))

Theories behind the manipulation of visualizations include but are not limited to - cognitive fit ((Vessey & Galletta, 1991)) - visual perceptual approaches ((**baker2009?**)) - human information processing

As interactive visualizations take a larger role in information systems, designers must know what tasks, visual representations, and interaction techniques are available and how they work in concert to facilitate analytical reasoning. They must decided on the most effective visual representation without being able to estimate every user's ability to read and interpret the visualization (**boy2014?**). Tasks can be viewed either by the goal the user is trying to obtain or by the intent the user has

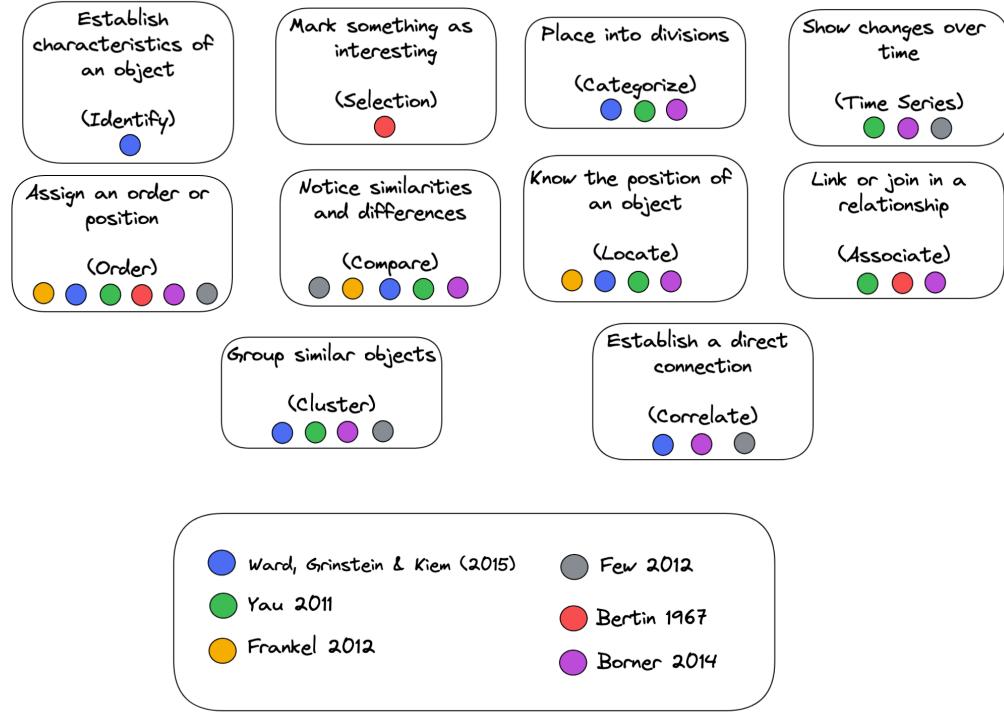


Figure 1.7: 10 Type Tasks with Sources

Interaction of Type Tasks facilitate data exploration leading to the generation of new insights. Interactions explicitly place humans in the loop where visualizations leverage the human perceptual system reducing the cognitive load required for data analysis ((Endert, Chang, North, & Zhou, 2015), (Sedig, Parsons, & Babanski, 2012)). More systems are using interactive visualizations, as opposed to static visualizations, which in turn requires a strong need to fully understand the effectiveness of interaction techniques (Saket, Srinivasan, Ragan, & Endert, 2017).



Figure 1.8: Interaction Types with Sources

(Schulz, Nocke, Heitzler, & Schumann, 2013) define two abstractions for the design of visualizations:

$$Data + Task = Visualization$$

$$Data + Visualization = Task$$

These abstractions demonstrate dependence between the data, visual representation, and the task. The more the user interacts with the visualization, they gain knowledge. The interactions allow user to be in control of their understanding by providing the flexibility to create new views that help him/her go beyond just the visual representation (kiem2008?). The field of

information visualization is continually adapting to changes with the big data revolution.

1.3.3 Advancement of graphical visualizations in dashboards

Visualizations have become more effective in recent years due to the pandemic and the Johns Hopkins University COVID-19 Dashboard (Dong E, 2022) **Dashboard**. We were glued to our computers, TVs, and phones for most of the world. As a result, we watched the dashboard change in real-time to adapt to the users' needs. In part to data growth and changing, the dashboard, as well as the visualizations, were needed in a condensed platform. The need to be concise and vastly informative can be a bit of a struggle when it comes to data visualizations. The human brain can only take in a set amount of data at a time from a table or a paragraph. The space of infographics has been a much better way of looking at data on a creative scale. While this may be a way of seeing the data in a friendly way, infographics need to include an interactive piece of data that many people would like to explore.

1.3.4 Applications

Two general applications in areas of visual inference have developed since the work of (Buja et al., 2009). These applications are actual methodology and methodology based on protocols. Actual methodology applications are used with alternative and corresponding null hypotheses, which perform visual inference tests to show many participants of different backgrounds with lineups.

1.4 Dashboard Design/ Human Perception (Visual Information)

1.4.1 History/Origin of Dashboard Design

A dashboard is a visual display of the most essential information needed to achieve one or more objectives, consolidated and arranged on a single screen so the information can be monitored at a glance(Few, 2006). A dashboard has particular characteristics:

- Achieve specific objectives
- Fits on a single computer screen
- Information can be displayed in multiple mediums (web browser or mobile device)
- Can be used to monitor information at a high level

While a dashboard can be extremely useful, it may be worth describing that a poorly designed dashboard will not be used. A dashboard should be concise, clear, and intuitive when displaying components in combination with a customized list of requirements of users.

Much of the work done within statistical research and dashboard design involves collaboration with other academic researchers. While this may be the best for the growth of the discipline, one will find that working with collaborators with non-STEM backgrounds.

A collaboration is when 2 or more entities work together to produce a desired and shared outcome. Interdisciplinary research with collaboration is a

pinnacle aspect of dashboards and statistical graphics to produce innovation and scientific knowledge.

The National Academies defines research collaboration as follows (NATIONAL ACADEMY OF SCIENCES, NATIONAL ACADEMY OF ENGINEERING, AND INSTITUTE OF MEDICINE OF THE NATIONAL ACADEMIES, 2004):

Interdisciplinary Research (IDR) is a mode of research by teams or individuals that integrates information, data, techniques, tools, perspectives, concepts, and/or theories from two or more disciplines or bodies of specialized knowledge to advance fundamental understanding or to solve problems whose solutions are beyond the scope of a single discipline or field of research practice.

1.4.2 Section 8: History/Origin of Human Perception of Statistical Graphics

Human perception plays a direct role in the area of visualization. The importance of human perception importance was cited by the NSF panel on graphics and image processing that proposed “scientific visualization” (NSF?).

Data Analysis tasks closely resemble the cognitive process known as sense-making. Tukey and Wilk (Tukey & Wilk, 1966) highlight the role of cognitive processes in their initial descriptions of EDA.

The primary general intent of data analysis is simply to seek through a body of data for exciting relationships and informa-

tion and to exhibit the results in such a way as to make them recognizable to the data analyzer.

Mallows and Walley (Mallows & Walley, 1980) list psychology as one of four areas likely to support a theory of analysis. Data analyses rely on the mind's ability to learn, analyze, and understand. Assigning meaning is not a statistical or computational step but a cognitive one. Each data analysis is part of a more extensive cognitive process.

Untrained analysts can and do “analyze” data with only their natural mental abilities - The mind performs its data analysis-like process to create detailed understandings of reality from bits of sensory input.

1.4.2.1 Schemas and Sensemaking

Cowan (Cowan, 2001) suggested that the average person can only hold two to six pieces of information in their attention. People can develop detailed understandings of reality, which is infinitely complex.

Cognitive structures consists of mental models and their relationships ((Rumelhart & Ortony, 1976), (Carley & Palmquist, 1992), (**jonassen1996?**)). Mental models have been studied under several different names, which are displayed in the following image of Cognitive Structures:

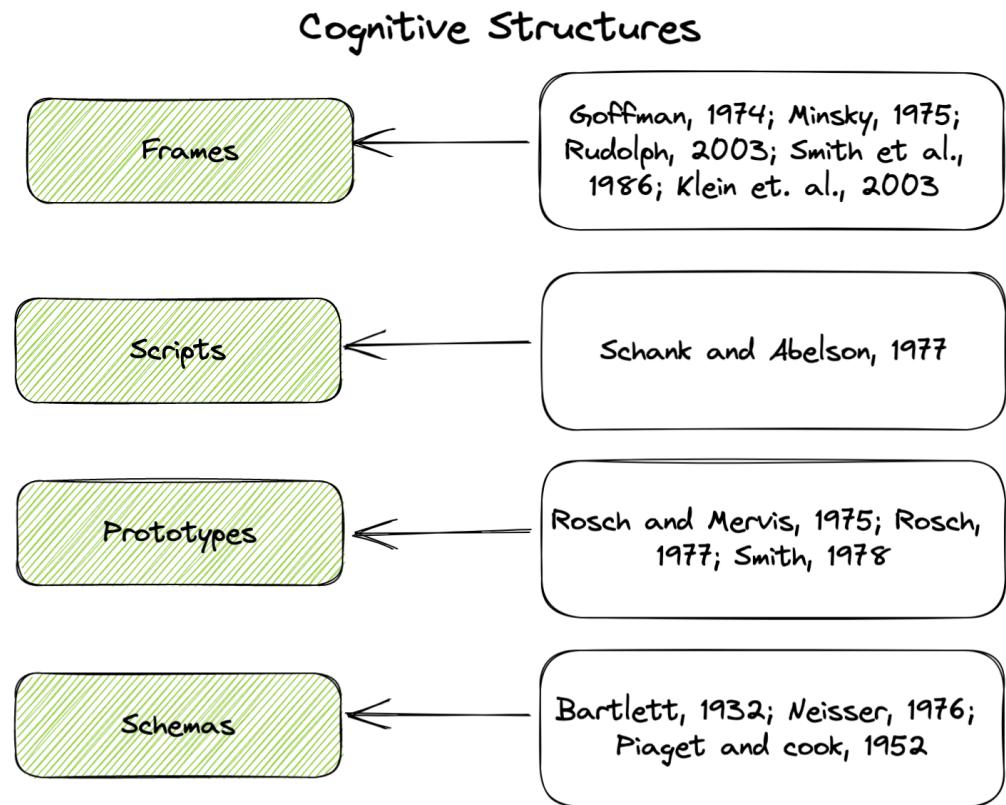


Figure 1.9: Cognitive Structures Citation Timeline

A schema is a mental model containing a breadth of information about a specific object or concept. Schemas are organized into semantic networks based on their relationships to other schemas (Wertheimer, 1938), (Rumelhart & Ortony, 1976). This arrangement helps the brain process its experiences instead of storing every sensory observation, the brain only needs to maintain its schemas, which are sufficient summaries of all previous observations. Some “memories” may even be complete recreations built with a schema (Bartlett & Remembering, 1932), (Klein, Phillips, Rall, & Peluso, 2007)

1.4.3 Section 9: Dashboard Design in regards to Human Perception

Data Scientists and Statisticians have produced more graphics since the pandemic's start. The reasons someone will create a graphic or dashboard may include but are not limited to understanding raw data structures to analyze model assumptions and present predictions along with displaying key performance metrics of business logic. These goals help work to navigate and are best served by quick-and-dirty representations of the data, while highly polished graphics may be more useful in other situations. It is valuable and essential to convey data correctly, meaning that we need to understand how graphics are perceived on a dashboard on a general level. Previous research by Tukey focused on graphics as a tool for exploratory analysis. Tukey describes in Exploratory Data Analysis (Tukey & Wilk, 1966) that pictures are often used to display data in a more enhanced version than a table. Tukey outlines detailed the types of different graphics and in which situations to utilize these graphics. The article - "External cognition: how do graphical representations work?" by Scaife and Rogers (Scaife & Rogers, 1996) critique the disparate literature on graphical representations, focusing on four representative studies. In general, this will help in the psychology of the perceptual experience.

The visual reference model developed by Card et al. (S. K. Card & Shneiderman, 1999) describes and identifies the three phases of the visualization process.

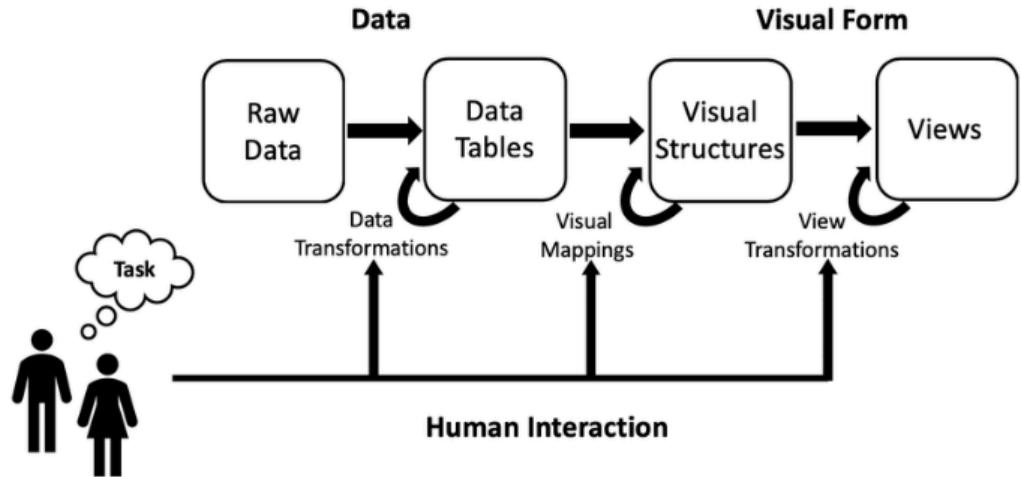


Figure 1.10: Visual Reference model by Card

While External Cognition describes the advances in graphical technology and how little had been done in the work of the cognitive framework of the discipline, the following citations by Ware attempt to develop the necessary guidelines that are useful for the work done by the perceptual experience.

Colin Ware “Information Visualization: Perception for Design” (Ware, 2004) there are four stages of visualization - The collection and storage of data itself - The preprocessing design to transform the data into something we can understand - The display hardware and the graphics algorithms produce the image on the screen. - The human perceptual and cognitive system

The overlapping understanding in the field, while Ware takes the process one step further not just to allow the end user to understand the outcomes but to curate the outcomes with a visual perception of the data that makes the cognitive load easier for the end-user.

A dashboard has much information related to tabular data from multiple sources. Design should be recorded a produced with content that will allow for

reproducibility. A dashboard should have some content related to interaction with the user. This interaction can be in multiple forms: toggling through the selection of variables to display uniformly to the use of interaction on the graph. Allowing a user to best understand what is being displayed in the graph.

The entire dashboard/Interface should have a human perception piece that is useful for the user to comprehend and use. If the user is visually overwhelmed, then the dashboard/interface is not practical.

- Identified the highly relevant from a dashboard design perspective (O'Donnell & David, 2000)
 1. Interaction and feedback are given by information systems
 2. Type of presentation format to be used
 3. Differences in the amount of information load.

Information load is essential, as dashboards must provide the right decision cues without overwhelming the user with excess information.

“A decision cue is a feature of something perceived that is used in the interpretation of perception (Choo, 2009),” where perception is an inferential process as objects in the environment can only be perceived indirectly through available information that has been sensed by the individual (Brunswik, 1952).

Visual complexity and information utility are required. Visual complexity refers to the ”degree of difficulty in providing a verbal description of an image.(Heaps & Handel, 1999), (Olivia, Mack, Shrestha, & Peepoer, 2004)

- Graphs are more suitable for spatial tasks (ex., for comparing a set of values) (Vessey & Galletta, 1991), (Umanath & Vessey, 1994); (Vessey, 1994)
- Graphs reduced the negative influence of information overload
- Graphs produced better correlation estimates and decreased time on a task (**schulzand?**; Booth & Siegler, 2006)
- Self-organizing maps and multidimensional scaling did not significantly outperform tabular representations (Huang et al., 2006)

The purposes of a dashboard:

1. Consistency
2. Monitoring
3. Planning
4. Communication

“Use interactive visual representations of abstract, non-physically based data to amplify cognition.” (Card, Mackinlay, & Shneiderman, 1999)

Visual perception involves two elements - the perceptual and conceptual gist. The perceptual gist refers to the process of the brain when it determines the image properties that provide the structural representation of a scene, like color and texture. The conceptual gist refers to the scene’s meaning, which is improved after the perceptual information is received. (Friedman, 1979); (Olivia, Mack, Shrestha, & Peepoer, 2004)

Visual complexity might increase with the quality and range of objects and with varying material and surface styles (Heylighen, 1997).

Repetitive and uniform patterns and existing knowledge of the objects in the scene reduce visual complexity (Olivia, Mack, Shrestha, & Peeper, 2004)

Guidelines for new generation dashboards refer to: 1. Aligning business processes with the latest information to provide business intelligence at all levels in the company. 2. Using intuitive and easy-to-digest visuals for delivering information to busy executives. 3. Sound Navigation (**ZiffDavisEnterprise2008?**)

If the guidelines on our visual information load can be related to statistical terminology, we can consider this the fisher information of visual information load. This load can be used as a metric for balancing the amount of information rather than overloading the consumer.

Visual Data Mining

For data mining to be effective, it is important to include humans in the data exploration process and combine the flexibility, creativity, and general knowledge of the computational power of today's computers. Visual data exploration aims at integrating humans in the data exploration process, applying their perceptual abilities to the large data sets available in today's computer systems.(Keim, 2002)

The visual data exploration process can be seen as a hypothesis generation process: - The visualizations of the data allow the user to gain insight into the data and come up with new hypotheses - Along with verification of the hypotheses

The main advantages of visual data exploration over automatic data mining techniques from statistics or machine learning are: - visual data exploration

can efficiently deal with highly inhomogeneous and noisy data. - visual data exploration is intuitive and requires no understanding of complex mathematical or statistical algorithms or parameters.

Visual Exploration Paradigm, also known as MGV (Massive Graph Visualizer), is an integrated visualization and exploration system for massive multidigraph navigation.(Abello & Korn, 2002). MGV usually follows a three-step process: - overview first - zoom and filter - details-on-demand

The user identifies interesting patterns and focuses on one or more of them. Note that visualization technology does not only provide the base visualization techniques for all three steps but also bridges the gap between the steps.

Visualization Technique Classification: - Standard 2D/3D displays such as bar charts and x-y plots - Geometrically transformed displays, such as landscapes and parallel coordinates, as used in a scalable framework - Icon-based displays such as needle icons and star icons as used in MGV - Dense pixel displays such as the recursive pattern and circle segments techniques and the graph sketches as used in MGV - Stacked displays, such as treemaps or dimensional stacking

Interaction and distortion techniques allow users to interact directly with the visualizations. - Projection as used in the Grand Tour System - Filtering as used in Polaris - Zooming as used in MGV and scalable framework - Linking and Brushing as used in Polaris and the scalable framework

Design Theory in Information System

The knowledge is distinguished as the fifth of five types of theory: 1. Analyzing & describing 2. Understanding 3. Predicting 4. Explaining and

predicting 5. Design and action

A definition of information systems that are suitable for our purposes concerns: "the effective design delivery use and impact of information technology in organizations and society (Avison, Fitzgerald, & DAWSON, n.d.).

The two paradigms characterize much of the research in the Information systems discipline: - *Behavioral Science Paradigm* - seeks to develop and verify theories that explain or predict human or organizational behavior (roots in natural science research methods). - *Data Science Paradigm* - seeks to extend the boundaries of human and organizational capabilities by creating new and innovative artifacts (roots in engineering and the sciences of the artificial).(Simon, 1996)

Technology and behavior are not dichotomous in an information system. They are inseparable (A. Lee, 2000). Information technology (IT) artifacts are broadly defined as: - constructs (vocabulary & symbols) - models (abstractions & representations) - methods (algorithms & practices) - instantiations (implemented & prototype systems)

These are concrete prescriptions that enable IT researchers and practitioners to understand and address the problems inherent in developing and successfully implementing information systems within organizations ((March & Smith, 1995); (Nunamaker, Dennis, Valacich, Vogel, & George, 1991)).

Interactive Visualization in Exploratory Data Analysis (EDA)

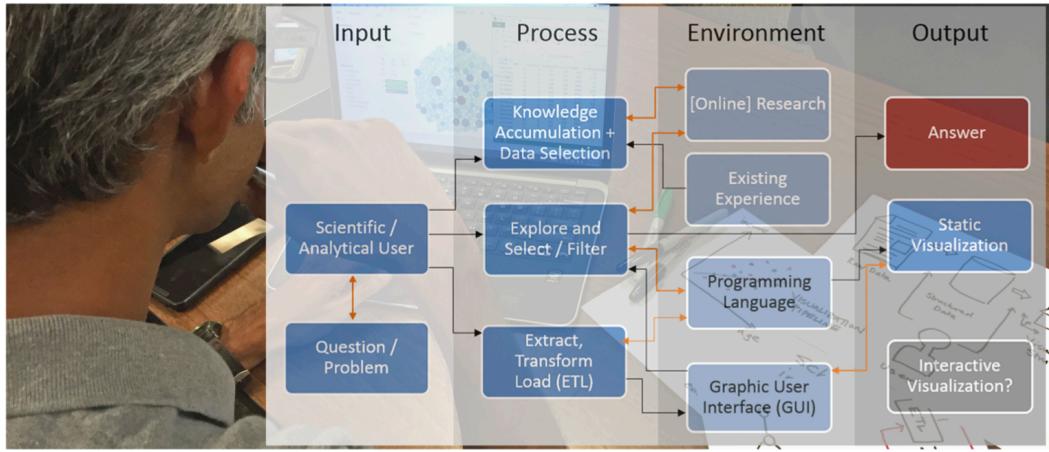


Figure 1. Elements of the exploratory data science workflow identified in this paper (background: analyst exploring a dataset).

Figure 1.11: Elements of EDA Workflow

Data Scientists and other analytic professionals often use interactive visualization in the dissemination phase at the end of a workflow, during which findings are communicated to a wider audience. Digital tools are critical to data science and analytics workflows, and current practice spans the following:

- *Data Analysis Tools* - R, Pandas and SAS
- *Data Warehousing Services*
- MySQL, MongoDB, or Amazon Redshift
- *Machine Learning Libraries* - scikit-learn o Apache MLlib

A common process typically consists of the following general stages:

1. *Discovery*: Formulating an exciting question and determining the data necessary to answer it.
2. *Acquisition*: Locating, organizing, and preparing data to be accessible to the chosen analysis environment.
3. *Exploration*: Investigating and analyzing the data set to collect insights and understand the data
4. *Modeling*: Building fitting and validating a model that can explain the data set and the observed phenomena
5. *Communication*: Disseminating the results to stakeholders in reports, presentations, and charts.

Static Visualization is commonly used in the communication phase of data

science workflows, and data scientists sometimes use them as part of the analysis. John Tukey’s EDA methods are known and well-vetted in the field currently. However, (Satyanarayan, Moritz, Wongsuphasawat, & Heer, 2016) began to address this by introducing a high-level grammar of graphics called “Vega-Lite,” which presents a set of standardized linguistic rules for producing interactive information visualizations using a concise JSON format for data to be represented by the grammar. Vega-Lite has been directly implemented in R via the `ggvis` package using the same - albeit slightly lower-level.

Contextual Design

(Wixon, Holtzblatt, & Knox, 1990) introduce “contextual design” as a systems development method in which the researcher partners with the user at the user’s place of work to “develop a shared understanding” of the user’s activities, and they define contextual inquiry as the first part of the broader process. Contextual inquiry is the data collection step of the field research element of the contextual design method, and it emphasizes four essential principles:

1. The context of the activity being performed by the user
2. The partnership between the researcher and the participant
3. The spoken verification that the investigator’s interpretation of the activity matches the user’s
4. The focus of the study is central to the approach taken by the interviewer

(Kandel, Paepcke, Hellerstein, & Heer, 2012) conducted what might be considered a contextual interview study similar to ours in that they ana-

lyzed data scientists' self-reported work processes. Kandal proposes three main archetypes that data scientists may be classed into the following:

- Hackers: who build processes chaining together multiple programming languages of different types (analytical, scripting, and database languages) and use visualization in various environments.
- Scripters: who perform most of their analysis in an analytical environment (e.g. R or Python) and perform the most complex statistical modeling of the types but who do not perform their ETL
- Application Users: who performed most or all of their work in an application such as Excel or SPSS and, like scripters, relied on others (namely, their organizations' IT departments) for ETL.

Dashboards can help understand and support many data types in essential business objectives. There are many different ways to label and utilize dashboards into different types.

Dashboards are cognitive tools that should be used to improve understanding of data, which should help people visually find relationships, trends, patterns, and outliers. Most importantly, dashboards should leverage people's visual cognitive capabilities.

EDA refers to methods and procedures for exploring the data space to learn about a data set. By analogy, exploratory modeling analysis (EMA) refers to methods and procedures for exploring the space of models which may be fit to a data set.

Interactive graphics are excellent for EDA, they are designed for exploring rather than presenting information (and more) and can be obtained by directly

querying the graphic (Unwin, Volinsky, & Winkler, 2003).

- PCPs enable the display of multi-dimensional data in two-dimensional space.
- There must be some loss of information, but this can be partly counteracted by varying the order of the axes.
- Interactivity is valuable for reordering the axes flexibly and fast.
- Interaction is valuable for dealing with the dense mass of lines produced by large data sets.

Being able to select subgroups of cases, highlight the selected lines, and switch between different subgroups all assist in interpreting the otherwise intricate displays which arise.

1.5 Conclusion

Chapter 2

Chapter Paper on Rural Shrink Smart Manuscript submitted to Journal of Data Science Special Issue

2.1 Abstract

Many small and rural places are shrinking. Interactive dashboards are the most common use cases for data visualization and context for exploratory data tools. In our paper, we will explore the specific scope of how dashboards are used in small and rural area to empower novice analysts to make data-driven decisions. Our framework will suggest a number of research directions to better support small and rural places from shrinking using an interactive dashboard design, implementation and use for the every day analyst.

2.2 Introduction

As the amount of data has increased in nearly every facet of life, the need to make sense of that data in an approachable, accessible form has become ever more important. As a result, many companies and organizations use interactive dashboards to present these data in a more useful and visually appealing form (Sarikaya, Correll, Bartram, Tory, & Fisher, 2019).

In many cases, dashboards support viewers' information processing, helping to make sense of complex data, navigate through a dataset, and supporting decision making based on the data.

Dashboards are often used, as with the car display of the same name, to provide summary information about many separate attributes of a common entity. One glance at a car's dashboard will tell you the speed, RPM, engine temperature, amount of gas in the tank; more importantly, however, the goal is not for the user to remember all of these characteristics, but to assess whether any of these quantities is outside of the expected range. Similarly, interactive dashboards for data are often used to display many different attributes and performance metrics which are of importance for stakeholders.

In this paper, we discuss the process of designing a dashboard to present publicly available government data to stakeholders in small Iowa towns to facilitate decision making and objective comparison with other similarly-situated towns.

Some communities continue to thrive as they lose population because they adapt, maintaining quality of life and community services for residents while investing in the future. This process, *smart shrinkage*, is important for rural areas who have experienced shrinking populations for decades. As small rural towns do not have access to data scientists or even the ability to easily leverage data collected locally to support decisions, our research team will provide communities with data about services in small town Iowa in order to assist with developing strategies to improve quality of life for their residents amid shrinking populations (Rural Shrink Smart Team, 2022). We hope to allow towns to explore their own data and compare to other similar towns, centering

decision-making on data in the context of small-town Iowa life.

2.3 Data Description

The Smart and Connected Community (SCC) dashboard data are primarily assembled from [data.iowa.gov](#) (State of Iowa, 2020), with some additional datasets assembled from federal and private sources. Most of these data sets are collected at a town/city or county spatial resolution, requiring us to carefully join data to ensure that these differences are respected while collating relevant information at the city level. In addition to the more commonly available statistics derived from e.g. the census and American Community Survey, [data.iowa.gov](#) contains several unique data sets, including local liquor sales, school building locations, town budgets and expenditures, hospital beds, Medicaid reimbursements, and other details that may provide information about local quality of life.

Data available on Iowa's data portal were augmented in some cases with higher-quality data sets in cases where the Iowa data were out of date or insufficiently accurate. Data collected from ELSI (National Center for Education Statistics, 2020) from <https://nces.ed.gov> were used to show the distance to any private or public school. The National Center for Education Statistics (NCES) is the primary federal entity for collecting and analyzing data related to education (Zarecor, Peters, & Hamideh, 2021).

Data collected from the Index of Relative Rurality (IRR) (USDA - ERS, 2020a) were used in the SCC dashboard to help classify the towns. The Index of Relative Rurality (IRR) is a continuous, threshold-free, and unit-free measure of rurality. It is an alternative to the traditional discrete threshold-based

classifications. The IRR ranges between 0 (low level of rurality, i.e., urban) and 1 (most rural). Four steps are involved in its design:

1. Identifying the dimensions of rurality: population size, density, remoteness, and built-up area.
2. Selecting measurable variables to adequately represent each dimension:
 - Size: logarithm of population size
 - Density: logarithm of population density.
 - Remoteness: network distance.
 - Built-up area: urban area (as defined by the US Census Bureau) as a percentage of total land area.
3. Re-scaling the variables onto bounded scales that range from 0 to 1.
4. Selecting a link function: unweighted average of the four re-scaled variable.

Data collected from Rural Urban Commuting Area Codes (USDA - ERS, 2020b) were used to help identify towns with commuting behaviors in our rural areas. The rural-urban commuting area (RUCA) codes classify U.S. census tracts using measures of population density, urbanization, and daily commuting. This data is on a zip code-level that will help identify those communities that commute to more urban areas. The most recent RUCA codes are based on data from the 2010 decennial census and the 2006-10 American Community Survey. The classification contains two levels. Whole numbers (1-10) delineate metropolitan, micropolitan, small town, and rural commuting areas based on the size and direction of the primary (largest) commuting flows.

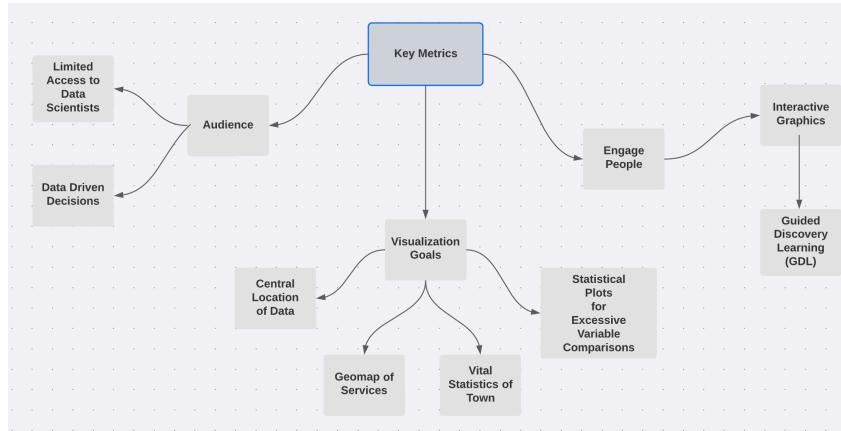


Figure 2.1: Diagram of considerations for our dashboard design process.

One of the interesting features of this assembled data set is that missing data can be missing for multiple reasons: not all state data is complete, but data about certain services may also be missing because towns do not offer that service. Thus, in addition to the usual challenges of working with real-world data that is “messy” in a variety of ways, we also have to contend with missing data that is missing due to the size of the community or the lack of services. This makes both visualization and statistical analysis more complicated (and more interesting).

2.4 Dashboard Design Considerations

One problem we identified early in the process of assessing smart-shrinkage strategies in small towns is that these towns do not have the resources to make data-driven decisions. Typically, small towns in Iowa are managed by at most a few part-time employees or volunteers. In some cases, essential management functions of the town are paid, but the municipalities we are interested in do not have sufficient funding to hire professionals to gather and analyze data.

As part of a wider project investigating the strategies towns use to main-

tain quality of life amid shrinking population, our research team provides communities with data about their own town, but also comparable towns across the state which may have a different approach to city services. In combination with other engagement strategies that are more qualitative, we hope to use this interactive dashboard approach to assist small Iowa cities with generalizing and developing strategies to improve or maintain quality of life amid shrinking populations.

One factor at the forefront of our visualization design is the importance of reducing the cognitive demands on viewers: we have assembled an incredible amount of data, and it is easy for even statisticians who deal with much larger datasets to get lost in the details of this data. At the same time, we want to invite viewers to engage with the data - to imagine, to draw comparisons, to generalize across towns, and to integrate outside information into the conclusions drawn based on the data we present. This invitation to engage with the data is similar to the approach advocated in Guided Discovery Learning, a framework leverages hints, feedback, and other helpful information to guide users in interactive exploration (DeDonno, 2016).

We expect that users will be interested in “sets” of variables from the wider dataset, which we assembled based on quality of life factors in the Iowa Small Town Poll (Peters, 2019). For instance, users might be interested in medical and social services available to residents, such as a local primary care clinic, nursing homes which are within driving distance, and the distance to the nearest emergency room; these factors might be explored separately from variables describing the services provided directly by city government, such as parks and recreation expenditures, snow removal services, and the distance to

the closest fire station.

As a consequence of this massively multivariate structure, we very quickly focused on the use of parallel coordinate plots; other alternatives, such as tours (Wickham, Cook, Hofmann, & Buja, 2011), require much more sustained attention to interactive plots as well as a deeper understanding of projections in multidimensional space which we cannot assume our users will have. Introduced in the 1880s (d’Ocagne, 1885), parallel coordinate or parallel set plots feature a series of vertical axes representing different variables arranged horizontally, with lines connecting each observation. When representing categorical data, parallel set plots may show “blocks” of data instead of individual lines, and are useful for representing conditional relationships between adjacent variables (Bendix, Kosara, & Hauser, 2005); modifications of this design, such as common-angle plots (Hofmann & Vendettuoli, 2013), address the issues which arise due to line-width illusions VanderPlas & Hofmann (2015). Parallel coordinate plots have been generalized to allow for continuous data and additional summaries beyond individual data points, such as densities (Heinrich & Weiskopf, 2009). In this paper, we use the `ggpcp` package, which leverages the grammar-of-graphics framework introduced in Wickham (2016), allowing us to use not only parallel coordinate plots, but also to overlay other statistical summaries, such as boxplots or violin plots, to provide additional context about the marginal distributions of each variable in addition to allowing for exploration of the multivariate space.

We also anticipate that users will be interested in comparing their town to other, similar towns. We will discuss the different ways that this comparison strategy was implemented in each dashboard in the next section, which de-

scribes the evolution of the dashboard over time and accounting for feedback from users and other researchers on the wider project.

One final component of this project is that our dashboard is part of a wider effort to work with towns to understand the different strategies used to maintain resident quality of life amid shrinking populations. Thus, while the town leaders are our primary audience, we also are creating this applet for use in parallel with a team of other researchers: sociologists, economists, city planning specialists, and artists. These researchers opinions and feedback about the dashboard are also useful and important, as they regularly work with town leaders in different capacities and have an understanding of what factors are most important to them and what types of questions these leaders may have when faced with data and unfamiliar statistical visualizations.

Throughout the design process, we will assess our visualizations to determine which strategies for user interface and interactive graphics design are most useful to empower town leaders to make discoveries in publicly available data assembled with a focus on items that impact rural quality of life.

2.5 Guiding Design Principles

Research on dashboard creation and interactive visualization tends to be very task-specific and hard to apply to more generalized settings. That is, it is relatively easy to create a dashboard that works for a particular task, but it is hard to generalize from that process what will work for the next dashboard. With this in mind, we set out to clearly document our intentions at each stage of the design and evaluation process, with the goal of gathering some useful information about general dashboard design from the process of creating this

specific dashboard.

Thus, our initial set of dashboard design principles is as follows:

- The town leaders are the focus audience; thus, the town itself should be the central focus of the app.
- We should facilitate comparisons with other towns in order to allow the user to explore other potential solutions to offering services that enhance resident quality of life.
- We will present the user with peer comparisons in order to widen the scope of exploration beyond the initial set of obvious peers in the local region.
- We will implement feedback mechanisms that allow us to provide more detailed data and respond to feature requests to improve the dashboard design over time.

As with many dashboards, this project is under continuous development; while it makes for an unsatisfactory conclusion, we do not have a “final” dashboard design because the application will continue to evolve. However, we have some useful insights into the process of creating an application designed to invite users to explore a large and complex dataset that we believe to be a useful contribution to work in this area.

2.6 Dashboard Design Process

2.6.1 Dashboard Components

In this section, we discuss the philosophy behind the basic “building blocks” of the dashboard. This philosophy is present in all of the iterations of the dashboard that we present in this discussion, and we will evaluate the overall philosophy’s effectiveness in the conclusion.

The large set of publicly available data (primarily from [data.iowa.gov](#)) we have assembled is useful, but we must be careful with how we present this data because it would be easy to overwhelm the user with small details that mask the bigger picture. We select a small subset of towns (out of the 999 towns in Iowa) and a small subset of variables of interest to start with, and then allow the user to increase the complexity of the display in accordance with their interest. This avoids some of the pitfalls of dashboard design that can easily lead to user overload (Few, 2006).

Our primary objective is to provide users with a town-centric approach: their town is at the center of our application, and comparisons to other, similar towns are secondary. As a result, the next component of the dashboard is intended to provide a brief overview of the information we have about a specific town of interest. This design is based on research into visualization sensemaking (S. Lee et al., 2016), in that we allow users to explore outward from the familiar to the unknown. The map visuals were built using Open Source Routing Machine (OSRM) route functions (Luxen & Vetter, 2011) in R (R Core Team, 2022) to amplify the accuracy of the distances from necessary services in town-centric point. OSRM allows for finding the “As the Crow Flies” distance and time on the road for our vital services map, since OSRM technology is similar to Google maps.

When faced with the next component, a parallel coordinate plot (PCP), a novice user will be able to determine two basic components: Visual Object (textual objects and non-textual objects) and Frame (frame of content and frame of visual encoding).

Taken together, the app is a single page; the initial “solid ground” which the user explores from consists of maps showing the route from the center of town to necessary services, including the fire department, schools, post offices, and hospitals. In version 2, as shown in [Figure 2.3](#), the map portion is condensed, and more space is given to value boxes that show vital statistics about the town’s QoL and financial metrics. This relatively straightforward display is followed by a parallel coordinate plot that allows the user to see similar towns along dimensions such as economic indicators or population size.

2.6.2 Initial Draft

The initial design sketch and implementation are shown in [Figure 2.2](#).

Users’ towns are at the center of our application, and comparisons to other, similar towns are secondary. As it can be extremely difficult to predict which towns are optimal for comparison purposes (similar may involve population, region, economic indicators, sports rivalries, and any number of other variables), we allow users to modify a set of suggested comparison towns to indicate other towns of interest.

We implemented some suggested town comparisons using unsupervised clustering methods to help our towns make decisions that are informed in comparison to similar towns, for budget size, population size and location. We initially focused on determining the next five to ten similar towns, based

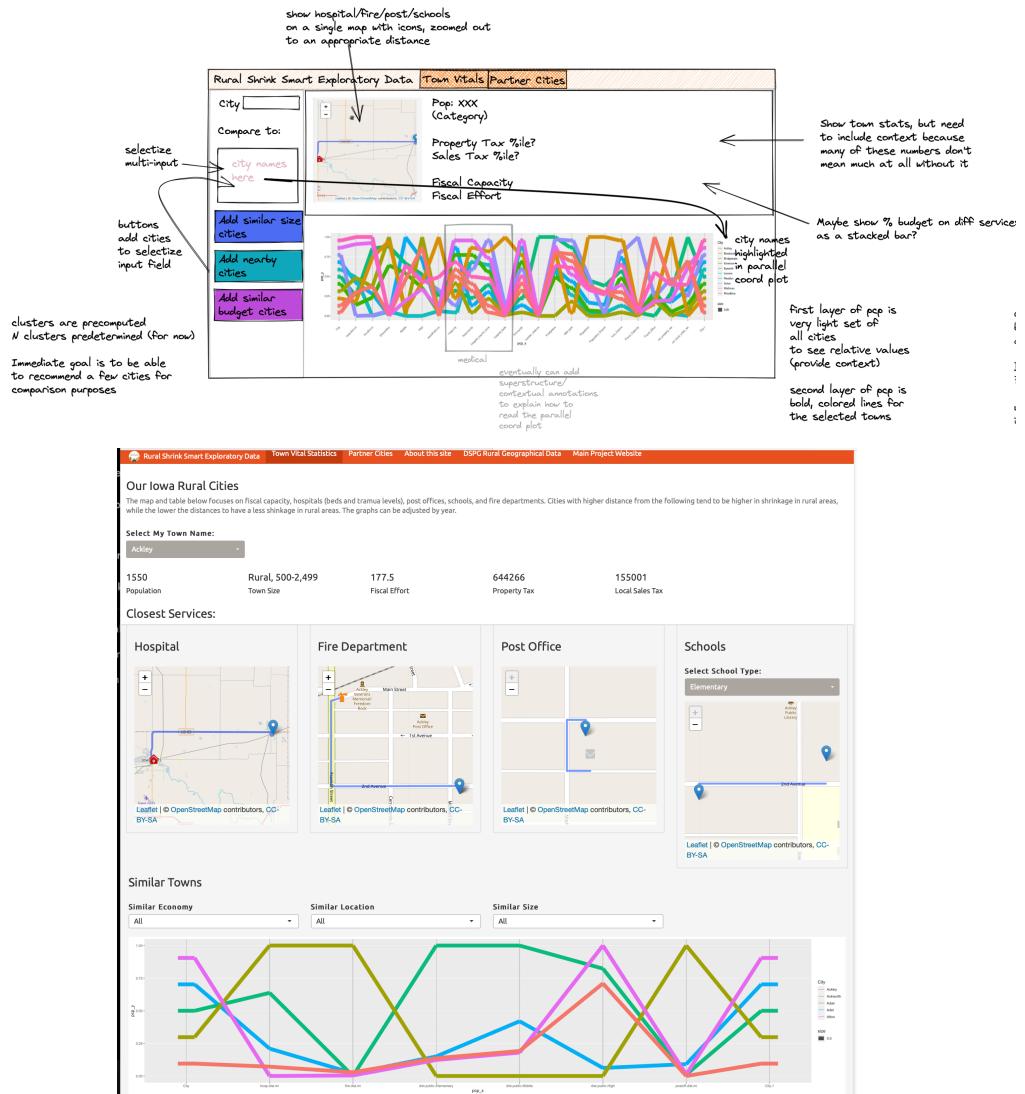


Figure 2.2: Initial dashboard design sketch (top) and implementation (bottom).

on distances to services. This feature became an important diagnostic for our data quality, as it became clear that towns which were grouped with big cities but which did not have a large population were so grouped because of missing data. Unfortunately, this clustering feature was not as useful to the application users, as they came to the dashboard with a pre-existing set of towns to compare to; our suggested comparisons were in the way.

The initial dashboard design featured several responsive maps showing the distance to the nearest hospital, fire department, post office, and school. These maps were ineffective for several reasons:

- Town residents already know this information (though it was useful for us as the dashboard designers, because we aren't nearly as familiar with the 900+ small towns in Iowa)
- We computed distance from services relative to the center of town - coordinates provided in the data from [data.iowa.gov](#). Generally speaking, the post office is at the center of town and the fire department is usually very close to the center of town; these two maps were useless. The school and hospital maps were less useless, but still did not provide particularly useful information to people already familiar with the town.
- It became clear that it might be more useful to show the comparison towns on a map (relative to the town of interest) so that users could compare geographical ratings for unfamiliar data to familiar data.

In addition, we received feedback on the parallel coordinate plot at the bottom of the app which was surprising: the viewers (in this case, other researchers on the team) were not as intimidated by the parallel coordinate plot

as we had expected. They did need some explanation of how to read the plot, and these hints need to be included in the dashboard, but they grasped the fundamental idea of the plot very quickly.

Our conclusion, based on this initial dashboard draft, was that we needed to restructure the application. Our attempt to show familiar information first to “build up” to the more unfamiliar structure of a parallel coordinate plot was not effective; there was too much clutter and not enough new information to draw users in.

2.6.3 Redesign

In the initial design, we included a map for each vital service, this initially created a lag for the users’ experience. As a result, we cached map directions from OSRM for each service in our database, which drastically reduced the response time for the user. Our initial design did not naturally focus the user’s eye on the most important parts of the dashboard; the redesign allowed for a cleaner flow from the top to the bottom.

In addition to the timing due to the map loading slowly, we added the vital statistics at the county level to allow for a more robust understanding of the town and it’s surroundings. The rurality index provided a better classification and the USDA sources allowed for the town to understand the impact of the closest major city due to commuting for work and shopping at larger stores not available within the town.

We also modified the parallel coordinate plots in several ways:

- Our x-axis had a large number of variables that we as researchers believed to be the most strongly associated with quality of life. However, there



Figure 2.3: A second iteration of the sketched design (top) and the implementation (bottom).

were still too many variables for users to successfully parse. We reduced the number of variables, focusing on variables that had the highest data quality, and we grouped these variables by quality of life factor (Peters, 2019).

- Originally, parallel coordinate bands were scaled based on the selected comparison towns. This had the effect of truncating the range of variables and over-emphasizing differences between selected towns relative to the overall range of each variable over all towns in our data set. We chose to show all towns in the data set in a very light α grey color to provide some information about the overall range of each variable. Unfortunately, even with the low- α value, this increased the visual complexity of the plot and confused users. Future iterations will likely make use of another aesthetic, such as boxplots or violin plots, to show the range of values for all towns, and then use lines only for towns that are selected by the user. This should strike a balance between visual complexity and representing the data accurately.
- We noticed that users did not make use of our suggested comparison towns, and so we removed that option in favor of allowing users to enter their own comparison towns directly. Users already had pre-determined towns they wanted to compare to, and our suggestions were just in the way.

While not all of these modifications were well received in our second round of user testing, the changes did incrementally move the dashboard display towards our goal of allowing users to explore the data and engage with it. We continued to be surprised with how well users reacted to the parallel coordi-

nate plots, which we initially thought might be too abstract for users unfamiliar with multivariate data displays, but the ability to compare towns across multiple dimensions and examine the similarities and differences between their approaches to different services seemed to be intuitive for users once they understood that each vertical axis was a different variable.

2.7 Discussion

Our dashboard design philosophy worked primarily to promote a town-centric approach application with comparisons to other similar towns being secondary. This approach created a way for the user to see their town information at the top of the page and to explore the PCP after reviewing their own town's essential statistics. The PCP in the lower part of the dashboard allowed for the user to see the plot and adjust to the fact that they could add more towns to the plot, providing an opportunity to explore the wider dataset from a base of familiar knowledge.

While we initially framed the design around guided discovery learning, the approach did not seem to suffice for our user base; instead, we found that users were more drawn to the unfamiliar from the start. We will likely leverage this in future iterations by using visual forms such as flower plots to draw the users in; even though these plots are not ideal for numerical display of data, the visual novelty and aesthetic appeal will provide some motivation to continue exploring and thinking about the data.

One factor that we have briefly considered and have seen hints of in our user feedback is that towns may not want to be compared negatively with other towns. While users have very definite ideas about which towns they would

like to compare to, we can always mask the town names and move back to comparisons based on town size and other factors (for instance, whether or not a town is the county seat is a factor that is important outside of population). Using this approach, we would label each town as “Town 1,” “Town 2,” and so on, which would eliminate some of the fears about negative comparisons, but would also remove some of the novelty of the data dashboard for our users and would prevent users from drawing on their own outside knowledge about each of the comparison towns.

We also recognize that we need to leverage the expertise of others in our research team: we are working with artists, researchers in architecture, economists, and sociologists; these researchers provide outside knowledge that we do not have and may be able to help us create insightful use-cases to showcase the app and teach towns how to use it. We can also leverage the app to connect users with our research team, providing additional value to those who use the applet and facilitating development of strategies for maintaining quality of life amid shrinking populations.

2.8 Future Work

One avenue we will explore in future iterations of the dashboard is to incorporate other dashboards generated by different groups within this project. This will create a wider field of information to explore: for instance, some of the additional work will focus on the 99 towns featured in the Iowa Small Town Poll; this will allow us to showcase survey-based measures of quality of life alongside the more objective measurements assembled in the dataset discussed in this paper. While at least one tab of this omni-dashboard will still

focus on wider EDA and discovery, we hope to incorporate other information as well to provide a more well-rounded data display encompassing most of the facets of this complex project.

We are also mindful of a distinction between “eye candy” and purpose-driven data visualization. While we have typically focused on the latter, there is certainly a place in our dashboard for the former as well. “Eye candy” visualization is intended to draw the viewer in and motivate them to explore; while these visualizations may not be particularly effective at communicating quantitative information, if they motivate the user to engage with the rest of the dashboard, they still serve a purpose. It is with this mindset that we intend to explore the use of flower plots - the artistic opportunities combined with the display of quantitative information (even in a form that isn’t optimal for quantitative comparisons) may be useful to engage viewers before transitioning to more useful data visualizations intended to provide accurate quantitative comparisons.

EDA can be a difficult for a variety of groups of people, novice users and experienced researchers. One of the more difficult components of this project has been clearly articulating the purposes of EDA to a diverse group of researchers unfamiliar with the concept. One of the most useful parts of this dashboard iteration process has been as an aid to data discovery: that is, the dashboard motivated us to find additional data sources and incorporate them into the project. Having conversations with other researchers about the EDA process helped to facilitate these conversations, as each discussion seemed to uncover additional data sources that someone remembered after looking at the dashboard. While this facet of the dashboard process may be difficult to study

formally, it would be an interesting avenue for investigation.

2.9 Conclusions

In this paper, we have documented the process of designing a dashboard for exploration and visualization of a large and complex data set assembled from many different sources. Our primary audience was leaders of small towns in Iowa, with a secondary audience of researchers in fields other than statistics collaborating on this project with us. Through the process of revising our dashboard, we found that the idea of guided discovery learning as implemented in our first version did not work as well as we had anticipated. It was more important to focus on allowing users to explore their questions about the dataset by facilitating user-driven comparisons and exploration, rather than attempting to anticipate user desires by providing comparison towns. In addition, we found that it would be more effective to draw users in with novel visual displays, as these seemed to attract more interest than providing known facts and an opportunity to explore outwards from an initial area of familiarity.

While it is hard to apply the findings from one fairly specific visualization project more widely, there is a lack of resources in this area that provide both design philosophies and actual analysis of user feedback in a qualitative sense. We have attempted to address this dearth of information by providing the design strategies, user feedback, and our planned and executed modifications, in the hopes that others facing the daunting challenge of designing a dashboard for EDA may learn something from our experiences.

Chapter 3

Chapter 2 Stuff

Chapter 4

Tables, Graphics, References, and Labels

4.1 Tables

By far the easiest way to present tables in your thesis is to store the contents of the table in a CSV or Excel file, then read that file in to your R Markdown document as a data frame. Then you can style the table with the `kable` function, or functions in the `kableExtra` pacakge.

In addition to the tables that can be automatically generated from a data frame in **R** that you saw in [R Markdown Basics] using the `kable` function, you can also create tables using *pandoc*. (More information is available at <http://pandoc.org/README.html#tables>.) This might be useful if you don't have values specifically stored in **R**, but you'd like to display them in table form. Below is an example. Pay careful attention to the alignment in the table and hyphens to create the rows and columns. Generally I don't recommend this approach of typing the table directly into your R Markdown document.

Table 4.1: Correlation of Inheritance Factors for Parents and Child

Factors	Correlation between Parents & Child	Inherited
Education	-0.49	Yes
Socio-Economic Status	0.28	Slight
Income	0.08	No
Family Size	0.18	Slight
Occupational Prestige	0.21	Slight

We can also create a link to the table by doing the following: Table 4.1. If you go back to [Loading and exploring data] and look at the `kable` table, we can create a reference to this max delays table too: Table ???. The addition of the `(\#tab:inher)` option to the end of the table caption allows us to then make a reference to Table `\@ref(tab:label)`. Note that this reference could appear anywhere throughout the document after the table has appeared.

4.2 Figures

If your thesis has a lot of figures, *R Markdown* might behave better for you than that other word processor. One perk is that it will automatically number the figures accordingly in each chapter. You'll also be able to create a label for each figure, add a caption, and then reference the figure in a way similar to what we saw with tables earlier. If you label your figures, you can move the figures around and *R Markdown* will automatically adjust the numbering for you. No need for you to remember! So that you don't have to get too far into LaTeX to do this, a couple **R** functions have been created for you to assist. You'll see their use below.

In the **R** chunk below, we will load in a picture stored as `uw.png` in our main directory. We then give it the caption of “UW logo,” the label of “`uwl-logo`,” and specify that this is a figure. Make note of the different **R** chunk options that are given in the R Markdown file (not shown in the knitted document).

```
include_graphics(path = "figure/unl.png")
```

Here is a reference to the UW logo: Figure 4.1. Note the use of the `fig:` code here. By naming the **R** chunk that contains the figure, we can then reference that figure later as done in the first sentence here. We can also specify the caption for the figure via the R chunk option `fig.cap`.



Figure 4.1: logo

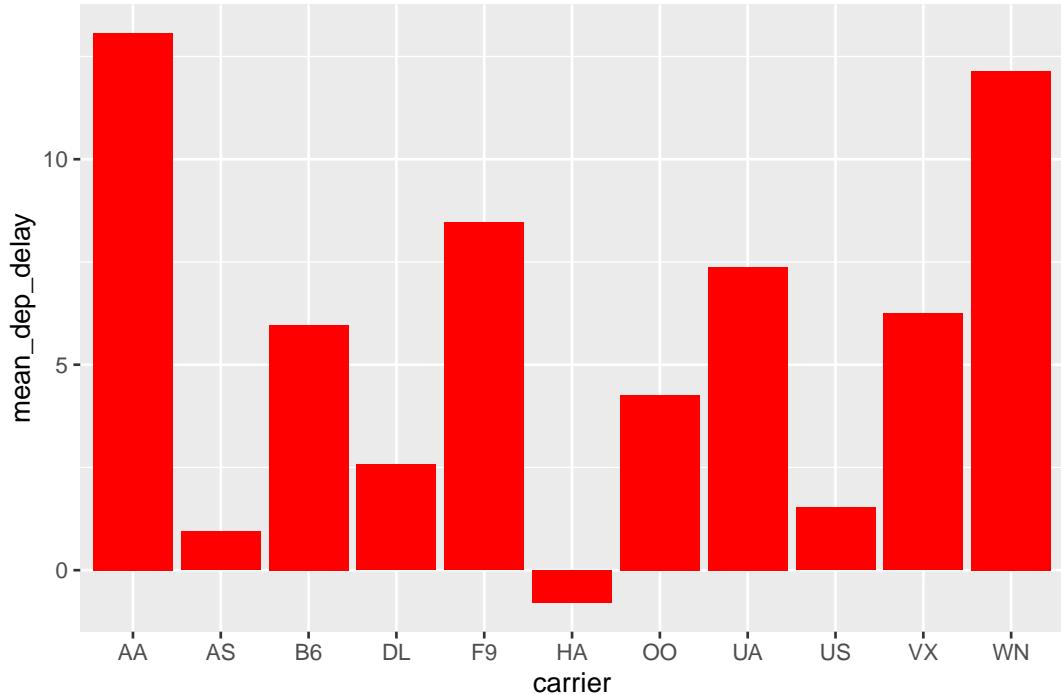


Figure 4.2: Mean Delays by Airline

Below we will investigate how to save the output of an **R** plot and label it in a way similar to that done above. Recall the **flights** dataset from Chapter 2. (Note that we've shown a different way to reference a section or chapter here.) We will next explore a bar graph with the mean flight departure delays by airline from Portland for 2014. Note also the use of the **scale** parameter which is discussed on the next page.

```
flights %>%
  group_by(carrier) %>%
  summarize(mean_dep_delay = mean(dep_delay)) %>%
  ggplot(aes(x = carrier, y = mean_dep_delay)) +
  geom_bar(position = "identity", stat = "identity",
           fill = "red")
```

Here is a reference to this image: Figure 4.2.

A table linking these carrier codes to airline names is available at <https://github.com/ismayc/pnwflights14/blob/master/data/airlines.csv>.

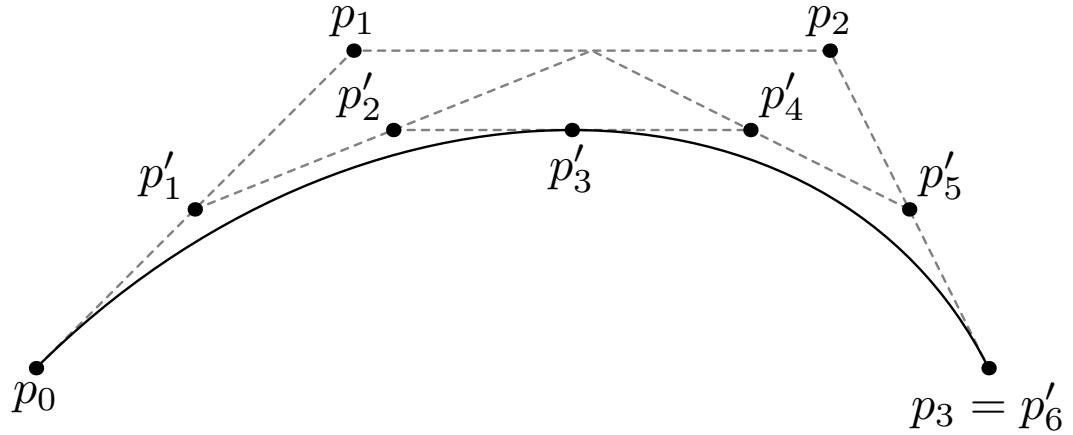


Figure 4.3: Subdiv. graph

Next, we will explore the use of the `out.extra` chunk option, which can be used to shrink or expand an image loaded from a file by specifying "`scale=`". Here we use the mathematical graph stored in the "subdivision.pdf" file.

Here is a reference to this image: Figure 4.3. Note that `echo=FALSE` is specified so that the **R** code is hidden in the document.

More Figure Stuff

Lastly, we will explore how to rotate and enlarge figures using the `out.extra` chunk option. (Currently this only works in the PDF version of the book.)

As another example, here is a reference: Figure 4.4.

4.3 Footnotes and Endnotes

You might want to footnote something.¹ The footnote will be in a smaller font and placed appropriately. Endnotes work in much the same way.

¹footnote text

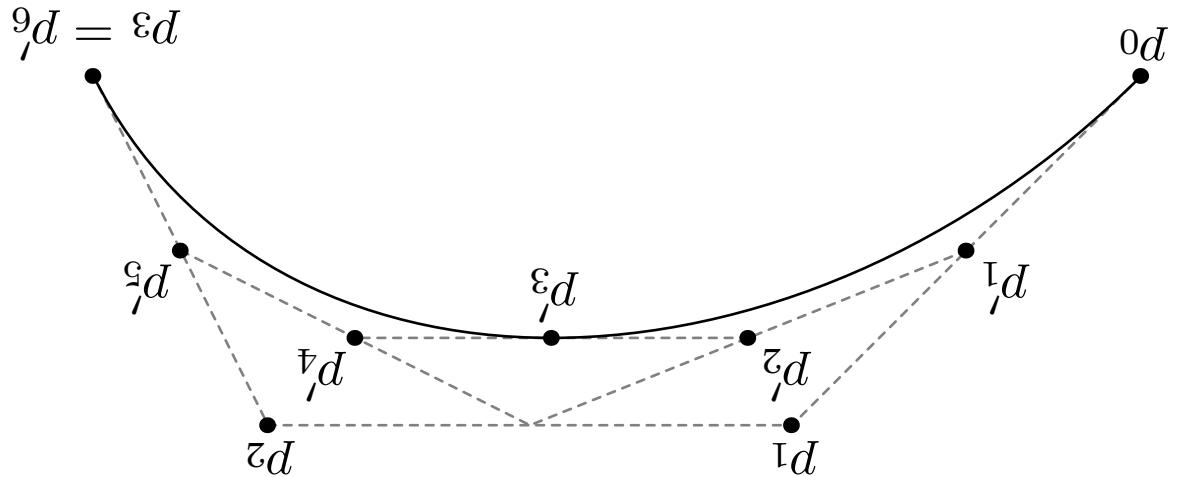


Figure 4.4: A Larger Figure, Flipped Upside Down

4.4 Cross-referencing chapters and sections

The [bookdown documentation](#) is an excellent source for learning how to cross-reference in a bookdown project such as a huskydown document. Here we only cover the most common uses for a typical thesis. If you want something more complex or fancy, please refer to the bookdown documentation and seek help from the developers of that package.

By default, all of your chapter and section headers will get an auto-generated ID label. For example, e.g., `# Chapter 1` will have an auto-generated ID `chapter-1`. Note that the ID label is all lower case, and has no spaces. If you have any kind of punctuation in your header, such as a colon (:), it will not appear in the ID label. Then in your text you can reference chapter one in your Rmd file like this: ‘as discussed in Chapter `\@ref(chapter-1)`,’ which will print as ‘as discussed in Chapter 1’

We strongly recommend that you manually assign ID labels to your chapter header to make it easy to cross-reference. For example, at the top of the Rmd file for this chapter, you can see:

```
# Tables, Graphics, References, and Labels {#ref-labels}
```

The {#ref-labels} part of this header is the ID label. It doesn't show in the output, but is there for us to use for easy cross-referencing, because it can be short, and we don't need to change it elsewhere our document when we update the chapter header. We can use this custom ID label in our Rmd document like this: 'as discussed in Chapter \@ref(ref-labels),' which will print as 'as discussed in Chapter 4.' If you need to show custom text instead of the chapter number, you use this syntax in your Rmd document: `see [my chapter about labels](#ref-labels) for more details` which will appear as 'see [my chapter about labels](#) for more details'

To cross-reference a specific section in the same chapter, we recommend adding a custom ID label to the section header, and using that to cross-reference. For example, earlier in this chapter we have a section on tables and in the Rmd file we see `## Tables {#tables}`. We can cross-reference that in the text like this 'as discussed in the section on [tables](#tables)' which will appear as 'as discussed in the above section on [tables](#)'

To cross-reference a section in a different chapter we can use the ID label from that section directly. For example, we can write in our Rmd document `as discussed in the section on [R code chunks](#r-chunks)` in Chapter \@ref(rmd-basics) which will appear as 'as discussed in the section on [R code chunks](#) in Chapter 2.'

If you prefer to cross-reference by the section number, we can use custom ID labels in our Rmd document. For example, to refer to a section in our first chapter, we can write in the Rmd document: `as discussed in section \@ref(r-chunks) in Chapter \@ref(rmd-basics)`. This will appear with

section and chapter numbers like so: as ‘as discussed in section ?? in Chapter 2.’

4.5 Bibliographies

Of course you will need to cite things, and you will probably accumulate an armful of sources. There are a variety of tools available for creating a bibliography database (stored with the .bib extension). In addition to BibTeX suggested below, you may want to consider using the free and easy-to-use tool called Zotero. Some Zotero documentation is at <http://libguides.reed.edu/citation/zotero>. In addition, a tutorial is available from Middlebury College at <http://sites.middlebury.edu/zoteromiddlebury/>.

R Markdown uses *pandoc* (<http://pandoc.org>) to build its bibliographies. One nice caveat of this is that you won’t have to do a second compile to load in references as standard LaTeX requires. To cite references in your thesis (after creating your bibliography database), place the reference name inside square brackets and precede it by the “at” symbol. For example, here’s a reference to a book about worrying: (Molina & Borkovec, 1994). This Molina1994 entry appears in a file called `thesis.bib` in the `bib` folder. This bibliography database file was created by a program called BibTeX. You can call this file something else if you like (look at the YAML header in the main .Rmd file) and, by default, is placed in the `bib` folder.

For more information about BibTeX and bibliographies, see (<http://web.reed.edu/cis/help/latex/index.html>)². There are three pages on this topic: *bibtex* (which talks about using BibTeX, at <http://web.reed.edu/>

²Reed College (2007)

[cis/help/latex/bibtex.html](http://web.reed.edu/cis/help/latex/bibtex.html)), *bibtexstyles* (about how to find and use the bibliography style that best suits your needs, at <http://web.reed.edu/cis/help/latex/bibtexstyles.html>) and *bibman* (which covers how to make and maintain a bibliography by hand, without BibTeX, at <http://web.reed.edu/cis/help/latex/bibman.html>). The last page will not be useful unless you have only a few sources.

If you look at the YAML header at the top of the main .Rmd file you can see that we can specify the style of the bibliography by referencing the appropriate csl file. You can download a variety of different style files at <https://www.zotero.org/styles>. Make sure to download the file into the csl folder.

Tips for Bibliographies

- Like with thesis formatting, the sooner you start compiling your bibliography for something as large as thesis, the better.
- The cite key (a citation's label) needs to be unique from the other entries.
- When you have more than one author or editor, you need to separate each author's name by the word "and" e.g. `Author = {Noble, Sam and Youngberg, Jessica},`.
- Bibliographies made using BibTeX (whether manually or using a manager) accept LaTeX markup, so you can italicize and add symbols as necessary.
- To force capitalization in an article title or where all lowercase is generally used, bracket the capital letter in curly braces.

4.6 Anything else?

If you'd like to see examples of other things in this template, please [contact us](#) (email bmarwick@uw.edu) with your suggestions. We love to see people using *R Markdown* for their theses, and are happy to help.

Conclusion

If we don't want Conclusion to have a chapter number next to it, we can add the `{-}` attribute.

More info

And here's some other random info: the first paragraph after a chapter title or section head *shouldn't be* indented, because indents are to tell the reader that you're starting a new paragraph. Since that's obvious after a chapter or section title, proper typesetting doesn't add an indent there.

Appendix A

The First Appendix

This first appendix includes all of the R chunks of code that were hidden throughout the document (using the `include = FALSE` chunk tag) to help with readability and/or setup.

In the main Rmd file

```
# This chunk ensures that the huskydown
# package is installed and loaded. This
# huskydown package includes the template
# files for the thesis.

if (!require(devtools)) install.packages("devtools",
  repos = "http://cran.rstudio.com")
# if(!require(huskydown))

# devtools::install_github(
#   'benmarwick/huskydown' )
# library(huskydown)

library(knitr)
library(palmerpenguins)
library(tidyverse)
```

In Chapter 4:

```
# This chunk ensures that the huskydown
# package is installed and loaded. This
# huskydown package includes the template
# files for the thesis and also two
# functions used for labeling and
# referencing

if (!require(devtools)) install.packages("devtools",
  repos = "http://cran.rstudio.com")

if (!require(dplyr)) install.packages("dplyr",
  repos = "http://cran.rstudio.com")

if (!require(ggplot2)) install.packages("ggplot2",
  repos = "http://cran.rstudio.com")

if (!require(bookdown)) install.packages("bookdown",
  repos = "http://cran.rstudio.com")

if (!require(huskydown)) {
  library(devtools)
  devtools::install_github("benmarwick/huskydown")
}

library(huskydown)

flights <- read.csv("data/flights.csv")
```

Appendix B

The Second Appendix, for Fun

Colophon

This document is set in **EB Garamond**, **Source Code Pro** and **Lato**. The body text is set at 11pt with *lmr*.

It was written in R Markdown and *LATEX*, and rendered into PDF using **huskydown** and **bookdown**.

This document was typeset using the XeTeX typesetting system, and the **University of Washington Thesis class** class created by Jim Fox. Under the hood, the **University of Washington Thesis LaTeX template** is used to ensure that documents conform precisely to submission standards. Other elements of the document formatting source code have been taken from the **Latex**, **Knitr**, and **RMarkdown templates for UC Berkeley's graduate thesis**, and **Dissertate: a LaTeX dissertation template to support the production and typesetting of a PhD dissertation at Harvard, Princeton, and NYU**

The source files for this thesis, along with all the data files, have been organised into an R package, `xxx`, which is available at <https://github.com/xxx/xxx>. A hard copy of the thesis can be found in the University of Washington library.

This version of the thesis was generated on 2022-12-03 09:13:03. The repository is currently at this commit:

The computational environment that was used to generate this version is as follows:

```
## - Session info -----
## setting value
##   version R version 4.1.0 (2021-05-18)
##   os       macOS Big Sur 10.16
##   system  x86_64, darwin17.0
##   ui       X11
##   language (EN)
##   collate en_US.UTF-8
##   ctype    en_US.UTF-8
##   tz       America/New_York
##   date     2022-12-03
##   pandoc  2.11.4 @ /Applications/RStudio.app/Contents/MacOS/pandoc/ (via rmarkdown)
##
## - Packages -----
##   package      * version date (UTC) lib source
##   assertthat    0.2.1   2019-03-21 [2] CRAN (R 4.1.0)
##   backports     1.4.1   2021-12-13 [1] CRAN (R 4.1.0)
##   bookdown      0.29.1  2022-09-18 [1] Github (rstudio/bookdown@4890be2)
##   broom        0.8.0   2022-04-13 [2] CRAN (R 4.1.2)
##   cachem       1.0.6   2021-08-19 [1] CRAN (R 4.1.0)
##   callr        3.7.0   2021-04-20 [2] CRAN (R 4.1.0)
##   cellranger    1.1.0   2016-07-27 [2] CRAN (R 4.1.0)
##   cli          3.4.1   2022-09-23 [1] CRAN (R 4.1.2)
##   colorspace    2.0-3   2022-02-21 [1] CRAN (R 4.1.2)
##   crayon       1.5.2   2022-09-29 [1] CRAN (R 4.1.2)
##   DBI          1.1.2   2021-12-20 [1] CRAN (R 4.1.0)
##   dbplyr       2.1.1   2021-04-06 [2] CRAN (R 4.1.0)
##   devtools      * 2.4.4   2022-07-20 [1] CRAN (R 4.1.2)
```

```
## digest          0.6.30  2022-10-18 [1] CRAN (R 4.1.2)
## dplyr          * 1.0.10  2022-09-01 [1] CRAN (R 4.1.2)
## ellipsis        0.3.2   2021-04-29 [2] CRAN (R 4.1.0)
## evaluate       0.18    2022-11-07 [1] CRAN (R 4.1.2)
## fansi           1.0.3   2022-03-24 [1] CRAN (R 4.1.2)
## farver          2.1.1   2022-07-06 [1] CRAN (R 4.1.2)
## fastmap         1.1.0   2021-01-25 [2] CRAN (R 4.1.0)
## forcats         * 0.5.1   2021-01-27 [1] CRAN (R 4.1.0)
## formatR          1.12    2022-03-31 [1] CRAN (R 4.1.2)
## fs               1.5.2   2021-12-08 [1] CRAN (R 4.1.0)
## generics         0.1.3   2022-07-05 [1] CRAN (R 4.1.2)
## ggplot2          * 3.3.6   2022-05-03 [2] CRAN (R 4.1.2)
## git2r            0.30.1  2022-03-16 [1] CRAN (R 4.1.2)
## glue              1.6.2   2022-02-24 [1] CRAN (R 4.1.2)
## gtable            0.3.0   2019-03-25 [2] CRAN (R 4.1.0)
## haven             2.5.0   2022-04-15 [2] CRAN (R 4.1.2)
## hms               1.1.1   2021-09-26 [1] CRAN (R 4.1.0)
## htmltools          0.5.3   2022-07-18 [1] CRAN (R 4.1.2)
## htmlwidgets        1.5.4   2021-09-08 [1] CRAN (R 4.1.0)
## httpuv             1.6.6   2022-09-08 [1] CRAN (R 4.1.2)
## httr               1.4.3   2022-05-04 [1] CRAN (R 4.1.2)
## huskydown         * 0.0.5   2022-01-16 [1] Github (benmarwick/huskydown@addb48e)
## jsonlite           1.8.3   2022-10-21 [1] CRAN (R 4.1.2)
## knitr              * 1.41    2022-11-18 [1] CRAN (R 4.1.2)
## labeling            0.4.2   2020-10-20 [2] CRAN (R 4.1.0)
## later              1.3.0   2021-08-18 [1] CRAN (R 4.1.0)
## lifecycle          1.0.3   2022-10-07 [1] CRAN (R 4.1.2)
## lubridate           1.8.0   2021-10-07 [1] CRAN (R 4.1.0)
## magrittr            2.0.3   2022-03-30 [1] CRAN (R 4.1.2)
## memoise             2.0.1   2021-11-26 [1] CRAN (R 4.1.0)
## mime                0.12    2021-09-28 [1] CRAN (R 4.1.0)
```

```
##  miniUI          0.1.1.1 2018-05-18 [1] CRAN (R 4.1.0)
##  modelr           0.1.8   2020-05-19 [2] CRAN (R 4.1.0)
##  munsell          0.5.0   2018-06-12 [2] CRAN (R 4.1.0)
##  palmerpenguins * 0.1.0   2020-07-23 [1] CRAN (R 4.1.0)
##  pillar            1.8.1   2022-08-19 [1] CRAN (R 4.1.2)
##  pkgbuild          1.3.1   2021-12-20 [1] CRAN (R 4.1.0)
##  pkgconfig          2.0.3   2019-09-22 [2] CRAN (R 4.1.0)
##  pkgload            1.3.0   2022-06-27 [1] CRAN (R 4.1.2)
##  prettyunits        1.1.1   2020-01-24 [2] CRAN (R 4.1.0)
##  processx           3.5.3   2022-03-25 [2] CRAN (R 4.1.2)
##  profvis            0.3.7   2020-11-02 [1] CRAN (R 4.1.0)
##  promises           1.2.0.1  2021-02-11 [1] CRAN (R 4.1.0)
##  ps                 1.7.0   2022-04-23 [2] CRAN (R 4.1.2)
##  purrr              * 0.3.5   2022-10-06 [1] CRAN (R 4.1.2)
##  R6                 2.5.1   2021-08-19 [1] CRAN (R 4.1.0)
##  Rcpp               1.0.9   2022-07-08 [1] CRAN (R 4.1.2)
##  readr              * 2.1.2   2022-01-30 [1] CRAN (R 4.1.2)
##  readxl             1.4.0   2022-03-28 [2] CRAN (R 4.1.2)
##  remotes            2.4.2   2021-11-30 [1] CRAN (R 4.1.0)
##  reprex              2.0.1   2021-08-05 [2] CRAN (R 4.1.0)
##  rlang               1.0.6   2022-09-24 [1] CRAN (R 4.1.2)
##  rmarkdown           2.16    2022-08-24 [1] CRAN (R 4.1.2)
##  rstudioapi         0.14    2022-08-22 [1] CRAN (R 4.1.2)
##  rvest               1.0.2   2021-10-16 [1] CRAN (R 4.1.0)
##  scales              1.2.0   2022-04-13 [2] CRAN (R 4.1.2)
##  sessioninfo        1.2.2   2021-12-06 [1] CRAN (R 4.1.0)
##  shiny               1.7.3   2022-10-25 [1] CRAN (R 4.1.2)
##  stringi             1.7.8   2022-07-11 [1] CRAN (R 4.1.0)
##  stringr              * 1.4.1   2022-08-20 [1] CRAN (R 4.1.2)
##  tibble              * 3.1.8   2022-07-22 [1] CRAN (R 4.1.2)
##  tidyverse            * 1.2.0   2022-02-01 [1] CRAN (R 4.1.2)
```

```
## tidyselect      1.2.0   2022-10-10 [1] CRAN (R 4.1.2)
## tidyverse       * 1.3.1   2021-04-15 [2] CRAN (R 4.1.0)
## tzdb            0.3.0   2022-03-28 [1] CRAN (R 4.1.0)
## urlchecker     1.0.1   2021-11-30 [1] CRAN (R 4.1.0)
## usethis        * 2.1.6   2022-05-25 [1] CRAN (R 4.1.2)
## utf8             1.2.2   2021-07-24 [2] CRAN (R 4.1.0)
## vctrs            0.5.1   2022-11-16 [1] CRAN (R 4.1.2)
## withr            2.5.0   2022-03-03 [1] CRAN (R 4.1.2)
## xfun             0.35    2022-11-16 [1] CRAN (R 4.1.2)
## xml2             1.3.3   2021-11-30 [1] CRAN (R 4.1.0)
## xtable           1.8-4   2019-04-21 [2] CRAN (R 4.1.0)
## yaml              2.3.6   2022-10-18 [1] CRAN (R 4.1.2)
##
## [1] /Users/dbradford4/Library/R/x86_64/4.1/library
## [2] /Library/Frameworks/R.framework/Versions/4.1/Resources/library
##
## -----
```

References

- Abello, J., & Korn, J. (2002). MGV: A system for visualizing massive multidigraphs. *IEEE Transactions on Visualization and Computer Graphics*, 8(1), 21–38. <http://doi.org/10.1109/2945.981849>
- Alsallakh, B., Gröller, M. E., Miksch, S., & Suntinger, M. (2011). Contingency wheel: Visual analysis of large contingency tables. In *EuroVA @ EuroVis*.
- Angel, E. (2000). *Interactive computer graphics : A top-down approach with OpenGL*. Boston, MA: Addison Wesley Longman.
- Angel, E. (2001a). *Batch-file computer graphics : A bottom-up approach with QuickTime*. Boston, MA: Wesley Addison Longman.
- Angel, E. (2001b). *Test second book by angel*. Boston, MA: Wesley Addison Longman.
- Ankerst, M., Berchtold, S., & Keim, D. A. (1998). Similarity clustering of dimensions for an enhanced visualization of multidimensional data. In *Proceedings IEEE symposium on information visualization (cat. no.98TB100258)* (pp. 52–60). <http://doi.org/10.1109/INFVIS.1998.729559>

- Avison, D., Fitzgerald, G., & DAWSON, C. (n.d.). Information systems development: Methodologies, techniques and tools, McGraw-hill.
- Bartlett, F. A., & Remembering, A. (1932). A study in experimental and social psychology. New York: Cambridge University Press.
- Bendix, F., Kosara, R., & Hauser, H. (2005). Parallel Sets: Visual Analysis of Categorical Data. In *2005 IEEE symposium on information visualization (INFOVIS'05)* (pp. 133–140). IEEE. <http://doi.org/10.1109/INFOVIS.2005.27>
- Beygelzimer, A., Perng, C.-S., & Ma, S. (2001). Fast ordering of large categorical datasets for better visualization. In *Proceedings of the seventh ACM SIGKDD international conference on knowledge discovery and data mining* (pp. 239–244).
- Booth, J. L., & Siegler, R. S. (2006). Developmental and individual differences in pure numerical estimation. *Developmental Psychology, 42*(1), 189.
- Bostock, M., Ogievetsky, V., & Heer, J. (2011). D³ data-driven documents. *IEEE Transactions on Visualization and Computer Graphics, 17*(12), 2301–2309.
- Brunswik, E. (1952). *The conceptual framework of psychology* (Vol. 1). University of Chicago Press.
- Buja, A., Cook, D., Hofmann, H., Lawrence, M., Lee, E.-K., Swayne, D. F., & Wickham, H. (2009). Statistical inference for exploratory data analysis and model diagnostics. *Philosophical Transactions of the Royal*

- Society A: Mathematical, Physical and Engineering Sciences*, 367(1906), 4361–4383.
- Card, S. K., Mackinlay, J. D., & Shneiderman, B. (1999). Using vision to think. *Readings in Information Visualization: Using Vision to Think*, 579–581.
- Carley, K., & Palmquist, M. (1992). Extracting, representing, and analyzing mental models. *Social Forces*, 70(3), 601–636.
- Choo, C. W. (2009). Information use and early warning effectiveness: Perspectives and prospects. *Journal of the American Society for Information Science and Technology*, 60(5), 1071–1082.
- Cleveland, W. S., & McGill, R. (1984). Graphical perception: Theory, experimentation, and application to the development of graphical methods. *Journal of the American Statistical Association*, 79(387), 531–554.
- Cowan, N. (2001). The magical number 4 in short-term memory: A reconsideration of mental storage capacity. *Behavioral and Brain Sciences*, 24(1), 87–114.
- d’Ocagne, M. (1885). Coordonnées parallèles et axiales : Méthode de transformation géométrique et procédé nouveau de calcul graphique déduits de la considération des coordonnées parallèles. *Gauthier-Villars*, 112. Retrieved from <https://archive.org/details/coordonnesparal00ocaggoog/page/n10>
- Day, R. H., & Stecher, E. J. (1991). Sine of an illusion. *Perception*, 20, 49–55.

- DeDonno, M. A. (2016). The influence of IQ on pure discovery and guided discovery learning of a complex real-world task. *Learning and Individual Differences*, 49, 11–16.
- Dong E, G. L., Du H. (2022). An interactive web-based dashboard to track COVID-19 in real time. [http://doi.org/10.1016/S1473-3099\(20\)30120-1](http://doi.org/10.1016/S1473-3099(20)30120-1)
- Endert, A., Chang, R., North, C., & Zhou, M. (2015). Semantic interaction: Coupling cognition and computation through usable interactive analytics. *IEEE Computer Graphics and Applications*, 35(4), 94–99.
- Few, S. (2006). *Information dashboard design: The effective visual communication of data*. Newton, MA: O'Reilly Media, Inc.
- Few, S. (2009). *Now you see it: Simple visualization techniques for quantitative analysis*. Analytics Press.
- Friedman, A. (1979). Framing pictures: The role of knowledge in automated encoding and memory for gist. *Journal of Experimental Psychology: General*, 108(3), 316.
- Friendly, M. (2014). Comment on “the generalized pairs plot.” *Journal of Computational and Graphical Statistics*, 23(1), 290–291.
- Fua, Y.-H., Ward, M. O., & Rundensteiner, E. A. (1999). Hierarchical parallel coordinates for exploration of large datasets. In *Proceedings visualization '99 (cat. no.99CB37067)* (pp. 43–508). <http://doi.org/10.1109/VISUAL.1999.809866>
- Guo, D. (2003). Coordinating computational and visual approaches for interactive feature selection and multivariate clustering. *Information*

- Visualization*, 2(4), 232–246.
- Heaps, C., & Handel, S. (1999). Similarity and features of natural textures. *Journal of Experimental Psychology: Human Perception and Performance*, 25(2), 299.
- Heinrich, J., & Weiskopf, D. (2009). Continuous Parallel Coordinates. *IEEE Transactions on Visualization and Computer Graphics*, 15(6), 1531–1538. <http://doi.org/10.1109/TVCG.2009.131>
- Heylighen, F. (1997). Publications on complex, evolving systems: A citation-based survey. *Complexity*, 2(5), 31–36.
- Hofmann, H., & Vendettuoli, M. (2013). Common Angle Plots as Perception-True Visualizations of Categorical Associations. *IEEE Transactions on Visualization and Computer Graphics*, 19(12), 2297–2305. <http://doi.org/10.1109/TVCG.2013.140>
- Huang, Z., Chen, H., Guo, F., Xu, J. J., Wu, S., & Chen, W.-H. (2006). Expertise visualization: An implementation and study based on cognitive fit theory. *Decision Support Systems*, 42(3), 1539–1557.
- Johansson, J., Ljung, P., Jern, M., & Cooper, M. (2005). Revealing structure within clustered parallel coordinates displays. In *IEEE symposium on information visualization, 2005. INFOVIS 2005.* (pp. 125–132). IEEE.
- Kandel, S., Paepcke, A., Hellerstein, J. M., & Heer, J. (2012). Enterprise data analysis and visualization: An interview study. *IEEE Transactions on Visualization and Computer Graphics*, 18(12), 2917–2926.

- Keim, D. A. (2002). Information visualization and visual data mining. *IEEE Transactions on Visualization and Computer Graphics*, 8(1), 1–8. <http://doi.org/10.1109/2945.981847>
- Klein, G., Phillips, J. K., Rall, E. L., & Peluso, D. A. (2007). A data-frame theory of sensemaking. In *Expertise out of context* (pp. 118–160). Psychology Press.
- Kolatch, E., & Weinstein, B. (2001). CatTrees: Dynamic visualization of categorical data using treemaps. Url: Http://www.Cs.Umd.Edu/-class/spring2001/cmsc838b/proje ct.Kolatch_Weinstein.
- Lee, A. (2000). Systems thinking, design science, and paradigms: Heeding three lessons from the past to resolve three dilemmas in the present to direct a trajectory for future research in the information systems field,“keynote address. In *Eleventh international conference on information management, taiwan*.
- Lee, S., Kim, S.-H., Hung, Y.-H., Lam, H., Kang, Y., & Yi, J. S. (2016). How do people make sense of unfamiliar visualizations?: A grounded model of novice’s information visualization sensemaking. *IEEE*, 22, 499–508.
- Luxen, D., & Vetter, C. (2011). Real-time routing with OpenStreetMap data. In *Proceedings of the 19th ACM SIGSPATIAL international conference on advances in geographic information systems* (pp. 513–516). New York, NY, USA: ACM. <http://doi.org/10.1145/2093973.2094062>
- Ma, S., & Hellerstein, J. L. (2001). Mining partially periodic event patterns with unknown periods. In *Proceedings 17th international conference on*

- data engineering* (pp. 205–214). IEEE.
- Mallows, C., & Walley, P. (1980). A theory of data analysis. *Proc. Amer. Statist. Assoc. Bus. Econ. Statist. Sec*, 8–14.
- March, S. T., & Smith, G. F. (1995). Design and natural science research on information technology. *Decision Support Systems*, 15(4), 251–266.
- Molina, S. T., & Borkovec, T. D. (1994). The Penn State worry questionnaire: Psychometric properties and associated characteristics. In G. C. L. Davey & F. Tallis (Eds.), *Worrying: Perspectives on theory, assessment and treatment* (pp. 265–283). New York: Wiley.
- NATIONAL ACADEMY OF SCIENCES, NATIONAL ACADEMY OF ENGINEERING, AND INSTITUTE OF MEDICINE OF THE NATIONAL ACADEMIES. (2004). FACILITATING INTERDISCIPLINARY RESEARCH. WASHINGTON, D.C.: THE NATIONAL ACADEMIES PRESS.
- National Center for Education Statistics. (2020). National center for education statistics. <https://nces.ed.gov/>.
- Nunamaker, J. F., Dennis, A. R., Valacich, J. S., Vogel, D., & George, J. F. (1991). Electronic meeting systems. *Communications of the ACM*, 34(7), 40–61.
- O'Donnell, E., & David, J. S. (2000). How information systems influence user decisions: A research framework and literature review. *International Journal of Accounting Information Systems*, 1(3), 178–203.

- Olivia, A., Mack, M. L., Shrestha, M., & Peeper, A. (2004). Identifying the perceptual dimensions of visual complexity of scenes. In *Proceedings of the annual meeting of the cognitive science society* (Vol. 26).
- Parsons, P., & Sedig, K. (2014). Adjustable properties of visual representations: Improving the quality of human-information interaction. *Journal of the Association for Information Science and Technology*, 65(3), 455–482.
- Peters, D. J. (2019). Community Resiliency in Declining Small Towns: Impact of Population Loss on Quality of Life over 20 Years. *Rural Sociology*, 84(4), 635–668. <http://doi.org/10.1111/ruso.12261>
- R Core Team. (2022). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing. Retrieved from <https://www.R-project.org/>
- Reed College. (2007). LaTeX your document. Retrieved from <http://web.reed.edu/cis/help/LaTeX/index.html>
- Rumelhart, D. E., & Ortony, A. (1976). The representation of knowledge in memory. In *Center for human information processing*.
- Rural Shrink Smart Team. (2022). Rural shrink smart. <https://ruralshrinksmart.org/>.
- S. K. Card, J. D. Mackinlay, & Schneiderman, B. (1999). Readings in information visualization: Using vision to think.
- Saket, B., Srinivasan, A., Ragan, E. D., & Endert, A. (2017). Evaluating interactive graphical encodings for data visualization. *IEEE Transactions*

- on Visualization and Computer Graphics*, 24(3), 1316–1330.
- Sarikaya, A., Correll, M., Bartram, L., Tory, M., & Fisher, D. (2019). What do we talk about when we talk about dashboards? *IEEE Transactions on Visualization and Computer Graphics*, 25(1), 682–692. <http://doi.org/10.1109/TVCG.2018.2864903>
- Satyanarayan, A., Moritz, D., Wongsuphasawat, K., & Heer, J. (2016). Vega-lite: A grammar of interactive graphics. *IEEE Transactions on Visualization and Computer Graphics*, 23(1), 341–350.
- Scaife, M., & Rogers, Y. (1996). External cognition: How do graphical representations work? *International Journal of Human-Computer Studies*, 45(2), 185–213.
- Schonlau, M. (2003). Visualizing categorical data arising in the health sciences using hammock plots. In *Proceedings of the section on statistical graphics, american statistical association*.
- Schulz, H.-J., Nocke, T., Heitzler, M., & Schumann, H. (2013). A design space of visualization tasks. *IEEE Transactions on Visualization and Computer Graphics*, 19(12), 2366–2375.
- Sedig, K., Parsons, P., & Babanski, A. (2012). Towards a characterization of interactivity in visual analytics. *J. Multim. Process. Technol.*, 3(1), 12–28.
- Simon, H. A. (1996). The sciences of the artificial 3rd ed. *MIT Press Cambridge*.
- State of Iowa. (2020). Iowa data portal. <https://data.iowa.gov>.

- Tukey, J. W., & Wilk, M. B. (1966). Data analysis and statistics: An expository overview. In *Proceedings of the november 7-10, 1966, fall joint computer conference* (pp. 695–709).
- Umanath, N. S., & Vessey, I. (1994). Multiattribute data presentation and human judgment: A cognitive fit perspective. *Decision Sciences*, 25(5-6), 795–824.
- Unwin, A., Volinsky, C., & Winkler, S. (2003). Parallel coordinates for exploratory modelling analysis. *Computational Statistics & Data Analysis*, 43, 553–564. [http://doi.org/10.1016/S0167-9473\(02\)00292-X](http://doi.org/10.1016/S0167-9473(02)00292-X)
- USDA - ERS. (2020a). Rural classifications. <https://www.ers.usda.gov/topics/rural-economy-population/rural-classifications/>.
- USDA - ERS. (2020b). Rural-urban commuting area codes. <https://www.ers.usda.gov/data-products/rural-urban-commuting-area-codes/>.
- Van Wijk, J. J. (2005). The value of visualization. In *VIS 05. IEEE visualization, 2005*. (pp. 79–86). IEEE.
- VanderPlas, S., & Hofmann, H. (2015). Signs of the sine illusion—why we need to care. *Journal of Computational and Graphical Statistics*, 24(4), 1170–1190. <http://doi.org/10.1080/10618600.2014.951547>
- Vessey, I. (1994). The effect of information presentation on decision making: A cost-benefit analysis. *Information & Management*, 27(2), 103–119.
- Vessey, I., & Galletta, D. (1991). Cognitive fit: An empirical study of information acquisition. *Information Systems Research*, 2(1), 63–84.

- Ware, C. (2004). *Information visualization: Perception for design*. San Francisco, CA: Morgan Kaufmann Publisher.
- Ware, C. (2012). *Information visualization: Perception for design*. Morgan Kaufmann.
- Wertheimer, M. (1938). Laws of organization in perceptual forms.
- Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis* (2nd ed.). Springer-Verlag New York. Retrieved from <https://ggplot2.tidyverse.org>
- Wickham, H., Cook, D., Hofmann, H., & Buja, A. (2011). tourr: An R Package for Exploring Multivariate Data with Projections. *Journal of Statistical Software, Articles*, 40(2), 1–18. <http://doi.org/10.18637/jss.v040.i02>
- Wickham, H., Hofmann, H., Wickham, C., & Cook, D. (2012). Glyph-maps for visually exploring temporal patterns in climate data and models. *Environmetrics*, 23(5), 382–393.
- Wixon, D., Holtzblatt, K., & Knox, S. (1990). Contextual design: An emergent view of system design. In *Proceedings of the SIGCHI conference on human factors in computing systems* (pp. 329–336).
- Yi, J. S., Kang, Y. ah, Stasko, J., & Jacko, J. A. (2007). Toward a deeper understanding of the role of interaction in information visualization. *IEEE Transactions on Visualization and Computer Graphics*, 13(6), 1224–1231.

Zarecor, K. E., Peters, D. J., & Hamideh, S. (2021). Rural smart shrinkage and perceptions of quality of life in the american midwest. In *Handbook of quality of life and sustainability* (pp. 395–415). Springer.