



Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

# Data Visualization

## Data 100: Principles and Techniques of Data Science

Sandrine Dudoit

Department of Statistics and Division of Biostatistics, UC Berkeley

Spring 2019



# Outline

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

## ① Motivation

## ② Principles of Data Visualization

2.1 Do We Really Need a Graph?

2.2 General Considerations

2.3 Graphical Perception

2.4 Bad Graphs

## ③ Survey of Data Visualization Techniques

3.1 One Quantitative Variable

3.2 Multiple Quantitative Variables

3.3 One Qualitative Variable

3.4 Multiple Qualitative Variables

3.5 Conditional Plots



# Data Visualization

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

*"One picture worth ten thousand words."*

Frederick R. Barnard, *Printer's Ink*, March 10th, 1927.



# An Oldie But Goodie

Data Visualization

Dudoit

Motivation

Principles of Data Visualization

Do We Really Need a Graph?

General Considerations

Graphical Perception  
Bad Graphs

Survey of Data Visualization Techniques

One Quantitative Variable

Multiple Quantitative Variables

One Qualitative Variable

Multiple Qualitative Variables

Conditional Plots

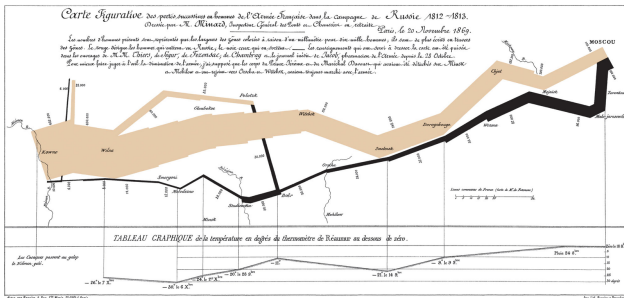


Figure 1: Minard's representation of Napoleon's 1812 Russian Campaign. This graph, made in 1861 by Charles Joseph Minard (1781–1870), is commonly regarded as one of the finest ever. It represents, in only two dimensions, the size of the troops, their location, their direction of movement, dates, and temperatures. [https://en.wikipedia.org/wiki/Charles\\_Joseph\\_Minard](https://en.wikipedia.org/wiki/Charles_Joseph_Minard).





# New But ...

Data Visualization

Dudoit

Motivation

Principles of Data

Visualization

Do We Really Need a Graph?

General Considerations

Graphical Perception

Bad Graphs

Survey of Data

Visualization

Techniques

One Quantitative Variable

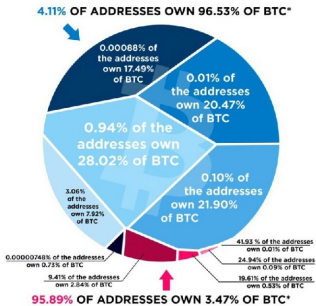
Multiple Quantitative Variables

One Qualitative Variable

Multiple Qualitative Variables

Conditional Plots

## The bitcoin Wealth Distribution



\* Data as of September 12th, 2013  
 Article and Sources:  
<https://howmuch.net/articles/bitcoin-wealth-distribution>  
<https://bitcoinprivacy.net/>



Figure 2: Bitcoin wealth distribution.

<http://viz.wtf/image/166329900475>.



# Data Visualization

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

One picture worth ten thousand words.

- Only if it is a good picture.
- We tend to be more demanding with text than with graphics.
- How long does it take to write/read one thousand words?  
At least the same effort should be put into making/viewing a graph.



# Learning Objectives

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Become a wise and effective “creator” / “maker” as well as “reader” / “viewer” of data visualization.
- Master general principles for data visualization and apply these when making your own graphs as well as when viewing others’.
- Produce the right graph for the matter at hand.
- Become aware of the variety of graphical techniques available for different types of data and purposes and understand their pros and cons.  
Go beyond histograms and pie charts!
- Think more carefully about each plot you create, consider the pros and cons of different choices, and try several different plots for a given dataset.



# Learning Objectives

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception  
Bad Graphs

Survey of  
Data  
Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Familiarize yourself with **software** for data visualization. Most of the examples in these slides are based on Python's `matplotlib` and `seaborn` libraries. However, as discussed in the first lecture, other languages such as R may be better suited for certain tasks.
- Focus on what type of plot to make rather than how to make it, i.e., **compose the plot conceptually before thinking of its software implementation** details. Concepts are general and long-lasting, while syntax is highly specific and ephemeral.
- **Avoid bad graphs!**



# Data Visualization

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception  
Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Data visualization is a **fundamental aspect of Data Science**.
- It is essential to “**look at data**” throughout the workflow, from exploratory data analysis (EDA) to model diagnostics and reporting the results of the inquiry.
- Visualization is valuable for **detecting the main features** (good or bad) of a dataset, **revealing patterns**, and **suggesting theories or further questions**.
- Visualization is also useful for **quality/assessment control** (QA/QC) and **detecting problems** with the data.
- An **effective plot can be good enough to answer the question on its own**. In some cases, it may even be the only appropriate type of answer.



# Data Visualization

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- An effective plot can also be sufficient to convince stakeholders of the findings from a full-blown statistical inference procedure.



# Data Visualization

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Although data visualization is ubiquitous and heavily relied upon, in research as well as in the media, typically **not much thought is put into creating or reading plots.**
  - ▶ **Creators** often rely on very limited subsets of plots and without proper consideration of their limitations.
  - ▶ **Readers** often passively absorb a message imposed on them by the graph, rather than reason and think critically about it.
- Very few Statistics, Computer Science (CS), or domain curricula offer courses in data visualization.
- Proper data visualization is **non trivial**. Entire courses could and should be devoted to data visualization, including discussions of **vision** and **perception** to guide the design of effective graphs.



# Do We Really Need a Graph?

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- When the data only comprise a handful of values, a **table** or a simple mention in **text** may be a more effective, i.e., accurate and simple, display.
- E.g. Percentage of popular vote for Trump and Clinton in 2016 presidential election:

Trump	46.1 %
Clinton	48.2 %





# Do We Really Need a Graph?

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

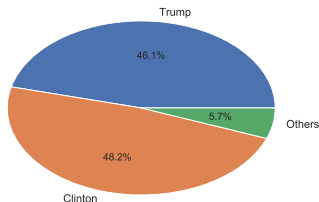
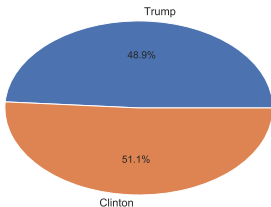


Figure 3: *US Election Results 2016*. Left: Pie chart of percentage of popular vote for Trump and Clinton. Right: Pie chart of percentage of popular vote for Trump, Clinton, and other candidates. **Why the different percentages on left and right?**



# From Tables to Graphs

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- When a table represents two or more variables, with more than a handful of values each, a **graph** may be more effective.
- Tables leave the interpretation to the viewer.
- Graphs provide a **summary** of the data and are more amenable to **comparisons**.
- Gelman et al. (2002). Lets Practice What We Preach: Turning Tables into Graphs. <http://www.stat.columbia.edu/~gelman/research/published/dodhia.pdf>.



# From Tables to Graphs

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

Profession	Frequency of recent citations	1996 total employed (1,000)	Relative frequency
Lawyers	8101	880	9.2
Economists	1201	148	8.1
Architects	1097	160	6.9
Physicians	3989	667	6.0
Statisticians	34	14	2.4
Psychologists	479	245	2.0
Dentists	165	137	1.2
Teachers (not university)	3938	4724	0.8
Engineers	934	1960	0.5
Accountants	628	1538	0.4
Computer programmers	91	561	0.2
Total	20,657	11,034	1.9

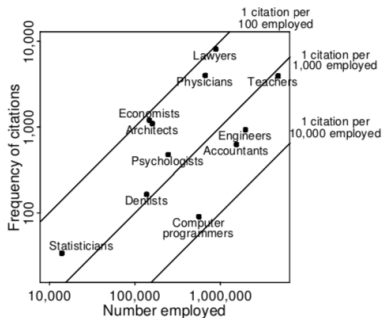


Figure 4: Turning tables into graphs (Gelman et al., 2002, Figure 2). Counts and rates of citations of various professions from the New York Times database. Graph: Log-log scale allows comparison across several orders of magnitude. Any 45° line indicates constant relative frequency. The relative positions of the different professions is clearer.



# More Oldies But Goodies

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots



Figure 5: *Album de Statistique Graphique (1881)*.

<https://www.davidrumsey.com/>.



# More Oldies But Goodies: Maps

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception  
Bad Graphs

Survey of  
Data  
Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

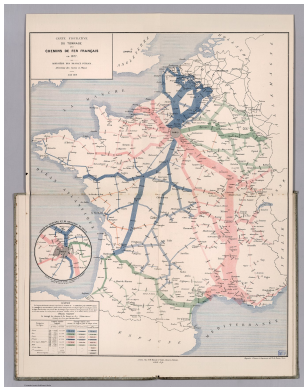


Figure 6: *Album de Statistique Graphique (1881)*. Train load (scaled by length of line) is represented by thickness of bands. How would you represent this data without a graph?



# More Oldies But Goodies: Graphical Timetables

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

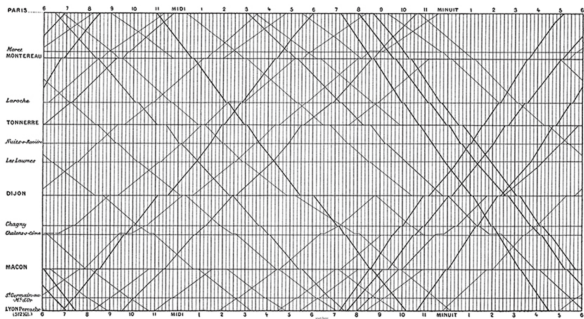


Figure 7: *Marey (1885). Train schedule Paris–Lyon, 1880s.*

[https://www.edwardtufte.com/bboard/q-and-a-fetch-msg?msg\\_id=0003zP](https://www.edwardtufte.com/bboard/q-and-a-fetch-msg?msg_id=0003zP). How would you represent this data without a graph?



# More Oldies But Goodies: Graphical Timetables

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

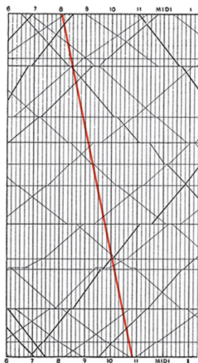


Figure 8: *Marey (1885). Train schedule Paris–Lyon with TGV, 1980s vs. 1880s.* The red line indicates the 1981 itinerary of the TGV, a new express train that cut the trip from Paris to Lyon to under three hours (vs. nine hours in the 1880s).



# More Oldies But Goodies: Graphical Timetables

Data Visualization

Dudoit

Motivation

Principles of Data Visualization

Do We Really Need a Graph?

General Considerations  
Graphical Perception  
Bad Graphs

Survey of Data Visualization Techniques

One Quantitative Variable  
Multiple Quantitative Variables  
One Qualitative Variable  
Multiple Qualitative Variables  
Conditional Plots

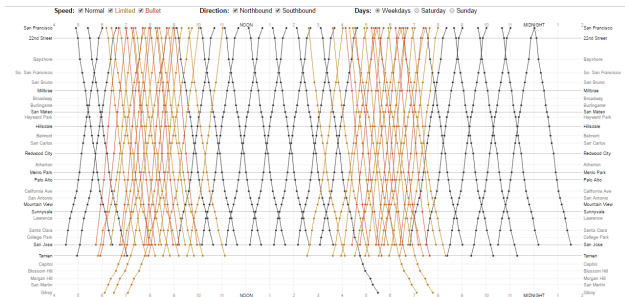


Figure 9: Train schedule SF–Gilroy, now.

<https://i.stack.imgur.com/qJ1hH>.





# More Oldies But Goodies: Graphical Timetables

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- In Marey (1885)'s Paris–Lyon **graphical train schedule** in the 1880s, time is represented on the x axis and the stations and distances between stations are represented on the y axis (Tufte, 2001).
- A train's **itinerary** is represented by a **line**.
- The **slope of the line** reflects the **speed** of the train: The more nearly vertical the line, the faster the train.
- The length of a stop at a station is indicated by the length of the horizontal line.
- The **intersection of two lines** locates the time and place that trains going in opposite directions pass each other.
- This type of graph, known as a **parallel coordinates plot**, is still used today and has many other applications.



# Caveats

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data  
Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Graphs should attempt to **summarize data in a simple, intuitive, and efficient manner, without distorting or losing important information.**
- However, **not all good graphs are simple.** As with text, plots conveying a lot of information (e.g., displaying multiple variables) require both a skillful creator and an educated reader.  
E.g. Minard's graph for Napoleon's Russia campaign, old graphical train schedules.
- There is **no "one-size-fits-all" graph**, i.e., **different types of graphs** should be used for different
  - ▶ **types of data**, e.g., quantitative, qualitative variables;
  - ▶ **purposes**, e.g., debugging code, EDA, reporting results;
  - ▶ **media**, e.g., print journal, projector.



# Caveats

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception  
Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Graphs typically **reduce the information** contained in the data.  
E.g. Histograms map  $n$  data points into  $B < n$  bins; boxplots map  $n$  data points into 5 summary statistics (+ possibly outliers).
- By **focusing on certain aspects** of the data or even **imposing structure** on data, graphs can also be **subjective**.  
E.g. Choosing which variables to plot, decisions regarding axes and scales, dendrogram representation of clusters<sup>1</sup>.
- As with text, the creator of the plot makes **editorial decisions** as to which data to display and which aspects of these data to show or emphasize.



# Caveats

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- The reader should **assess the relevance and reliability of the data** being displayed, as well as the **appropriateness of the graph**.
- **Software implicitly makes many decisions** for the creator of a plot, e.g., axes, scales, plotting symbols, color, ordering of data. **Experiment with different settings**.
- Graphs are rarely presented on their own. They should be **interpreted in context** of the text which they support. The reader should examine the **graph-text interface** and, in particular, whether the conclusions in the text are supported by the graph.

---

<sup>1</sup>A dendrogram is a graphical representation of hierarchical clustering results; for a given clustering of  $n$  objects, there are  $2^{n-1}$  possible dendrograms. The various choices made in hierarchical clustering as well as the dendrogram representation impose (vs. reveal) structure on the data.



# Statistical Inference

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Graphs are by definition **functions of the data**, i.e., **statistics**.
- Although not typically viewed this way, visualization can therefore be used as part of **statistical inference**.
- One can produce the **same types of plots for a sample and for a population**, in that sense, the plot for the sample can be viewed as an **estimator** of the plot for the population, i.e., the **parameter**.
- A pattern that we detect from plotting data for a sample can be used to **infer properties of the population** from which the sample was drawn. A formalized special case of such an approach is given by **linear regression**.



# General Considerations

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

In the process of creating a plot, you should consider the following issues.

- Determine the **purpose of the plot**.  
E.g. EDA, debugging code, comparing distributions, model diagnostics, summarizing results, reporting results.
- Formulate the **message**.
- Identify the **audience**.
- Identify the **display mode/medium** (e.g., journal, projector).
- Think about the **best type of graph** for the purpose, message, audience, and display mode.
- Aim for **efficient perception**: Speed, accuracy, and minimum cognitive load for understanding the message.



# General Considerations

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Apply **visual perception principles**.  
E.g. Angles and areas are harder to perceive/compare than lengths.
- Do not use more dimensions to represent the data than are in the data. This rules out pie charts and barplots.
- An important consideration when selecting a graphical technique is how easily it can be **extended** (e.g., to multiple variables) and how amenable it is to **comparing distributions**.
- Choose **graphical parameters** carefully: Aspect ratio, plotting symbols, line types, texture, axes, etc.



# General Considerations

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Choose **color palette** carefully. E.g. Be mindful of color blindness, use different color schemes for different types of data and messages (e.g., sequential, qualitative, and diverging).
- Provide **sufficient information** so that the plot can be interpreted properly.  
E.g. Title, axis parameters (i.e., label, tick marks), annotation, legend, caption, etc.  
In a document, number the figures and tables.
- Do not include irrelevant information, i.e., avoid “**chart junk**”.
- **Principle of “least surprise”**: If you defy expectations, people may get confused. Only defy expectations if it is very important.





# General Considerations

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data  
Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- **Experiment**, i.e., consider different types of plots and update the plots **iteratively**.
- Of course, always think about the **quality of the data** you plot.



# General Considerations

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One

Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- **Sample size.**
  - ▶ For small sample sizes, plot all of the data – Why loose information?
  - ▶ For larger samples sizes, plot relevant summaries of the data, that do not distort or loose important information in the data.
- **Variables to display/emphasize.** Depends on the purpose and message of the plot.
- **Type of variables.** Quantitative and qualitative variables call for different types of graphical summaries.
- **Pre-processing.** E.g. Transformation (e.g., log), dimensionality reduction, imputation.



# Graphical Perception

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Cleveland and McGill (1985): *“Graphical perception is the visual decoding of the quantitative and qualitative information encoded on graphs. Recent investigations have uncovered basic principles of human graphical perception that have important implications for the display of data.”*
- When we create a graph, we **encode** the data as **graphical attributes**.
- Possible **graphical attributes** are: Angles, areas, lengths, position on common aligned/unaligned scale, slopes, color properties.
- **Effective graphs** are those for which **attributes are most easily decoded**.



# Graphical Perception

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- There are **empirical laws for perception** that can be used to rank different types of graphical encodings.
- In general, such laws relate the **perceived (change in) intensity** in a physical stimulus to the **actual (change in) intensity**. This concerns stimuli to all senses, i.e., vision, hearing, taste, touch, and smell.



# Graphical Perception: Weber's Law

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- **Weber's Law** is an empirical relationship in psychophysics between the **initial intensity in a stimulus** ( $I$ ) and the **smallest perceivable difference** (a.k.a., just noticeable difference) in the stimulus intensity ( $\Delta I$ ):

$$\frac{\Delta I}{I} = k, \quad (1)$$

where  $k$  is a proportionality constant for a given type of stimulus <sup>2</sup>.

- In terms of length, this means we detect a 1 cm change in a 1 m length as easily as we detect a 10 m change in a 1 km length.
- Weber's Law appears to hold for many different graphical encodings.

---

<sup>2</sup>Law formulated and published by Gustav Theodor Fechner (1801–1887), a student of Ernst Heinrich Weber (1795–1878).



# Graphical Perception: Stevens' Law

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Stevens (1957) Law is an empirical relationship in psychophysics between the **intensity in a stimulus** and the **perceived magnitude of the sensation** created by the stimulus:

$$\psi(I) = Ci^\beta, \quad (2)$$

where  $I$  is the intensity or strength of the stimulus in physical units (energy, weight, pressure, mixture proportions, etc.),  $\psi(I)$  is the magnitude of the sensation,  $\beta$  is an exponent that depends on the type of stimulation or sensory modality, and  $C$  is a proportionality constant that depends on the units used.

- Examples of values for exponent,  $\beta$ 
  - Length: 0.9 – 1.1
  - Area: 0.6 – 0.9
  - Volume: 0.5 – 0.8



# Graphical Perception: Stevens' Law

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- For **lengths**, the relationship is almost **linear**, thus our perception is about right.
- However, according to this power law, our **perception of areas and volumes is conservative**, i.e., when values are represented as areas or volumes, we underestimate the large values relative to the small ones and overestimate the small ones relative to the large ones.
- E.g. Areas, with  $\beta = 0.7$ .  
Consider two areas of size 1 and 2, respectively.

$$\frac{\psi(2)}{\psi(1)} = \frac{2^{0.7}}{1^{0.7}} \approx 1.62.$$

Thus, we don't see the bigger area as twice as large.



# Graphical Perception: Stevens' Law

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

Now consider two areas of size  $1/2$  and  $1$ , respectively.

$$\frac{\psi(1/2)}{\psi(1)} = \frac{0.5^{0.7}}{1^{0.7}} \approx 0.62.$$

Thus, we don't see the smaller area as half as large.





# Graphical Perception: Stevens' Law

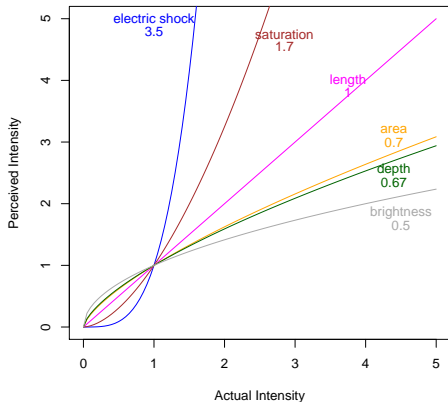


Figure 10: Graphical perception: Steven's Law. Stevens (1957) perceived sensory magnitude power law.



# Graphical Perception: Combining Weber's and Stevens' Laws

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Consider comparing the values  $x$  and  $x + w$ , using **length** ( $\beta = 1$ ) and **area** ( $\beta = 0.7$ ) encodings.
- For length, we perceive the relative value

$$\frac{x + w}{x} = 1 + \frac{w}{x}.$$

- For area, we perceive the relative value

$$\frac{(x + w)^{0.7}}{x^{0.7}} = \left(1 + \frac{w}{x}\right)^{0.7} \approx 1 + \frac{0.7w}{x}.$$

- Thus, we are **more likely to detect small differences using length encoding.**



# Graphical Perception

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Cleveland and McGill (1985) carried out an extensive study of graphical encodings to obtain a best to worst ranking.
- The encodings they examined include: position on a common aligned scale, position on a common unaligned scale, length, slope, angle, area, volume, color hue, brightness, and purity.
- One of their experiments consisted of
  - ▶ 7 graphical encodings,
  - ▶ 3 judgments per encoding,
  - ▶ 10 replications per subject,
  - ▶ 127 experimental subjects.

Assessment criterion:  $\text{error} = \|\text{perceived } p - \text{true } p\|$ ,  
where  $p$  denotes the ratio (in percentages) of the smaller to the larger magnitude.



# Graphical Perception

Data Visualization

Dudoit

Motivation

Principles of Data Visualization

Do We Really Need a Graph?

General Considerations

Graphical Perception

Bad Graphs

Survey of Data Visualization

Techniques

One Quantitative Variable

Multiple Quantitative Variables

One Qualitative Variable

Multiple Qualitative Variables

Conditional Plots

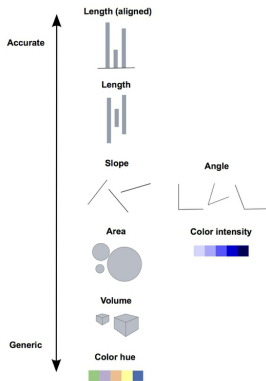


Figure 11: Graphical perception. Based on Table 1 in Cleveland and McGill (1985).

<http://paldhous.github.io/ucb/2016/dataviz/week2.html>.



# Bad Graphs

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

**Bad Graphs**

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- The literature is full of “bad graphs”, that, for instance, distort the data and are misleading, are too complicated, or are missing essential information.
- Karl Broman’s Top Ten Worst Graphs (including one of his own!): [https://www.biostat.wisc.edu/~kbroman/topten\\_worstgraphs/](https://www.biostat.wisc.edu/~kbroman/topten_worstgraphs/).
- Ross Ihaka’s Good and Bad Graphs: <https://www.stat.auckland.ac.nz/~ihaka/120/Lectures/lecture03.pdf>.
- Edward Tufte: [https://www.edwardtufte.com/bboard/q-and-a-fetch-msg?msg\\_id=00040Z](https://www.edwardtufte.com/bboard/q-and-a-fetch-msg?msg_id=00040Z).
- Junk Charts: [https://junkcharts.typepad.com/junk\\_charts/](https://junkcharts.typepad.com/junk_charts/).
- WTF Visualization: <http://viz.wtf>.



# Bad Graphs: Pie Charts

Data Visualization

Dudoit

Motivation

Principles of Data Visualization

Do We Really Need a Graph?

General Considerations

Graphical Perception

Bad Graphs

Survey of Data Visualization Techniques

One Quantitative Variable

Multiple Quantitative Variables

One Qualitative Variable

Multiple Qualitative Variables

Conditional Plots

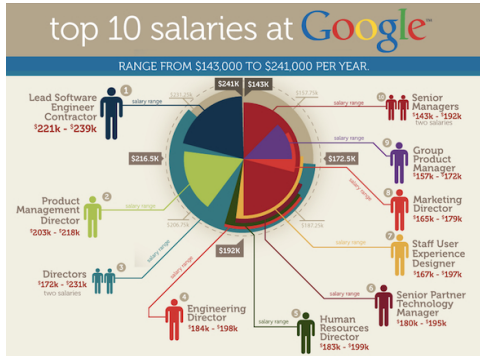


Figure 12: Top 10 Google salaries by job category: Pie chart.

[https://junkcharts.typepad.com/junk\\_charts/2011/10/the-massive-burden-of-pie-charts.html](https://junkcharts.typepad.com/junk_charts/2011/10/the-massive-burden-of-pie-charts.html)

What's the message? What do the angles represent? What's a better graph?



# Bad Graphs: Pie Charts

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

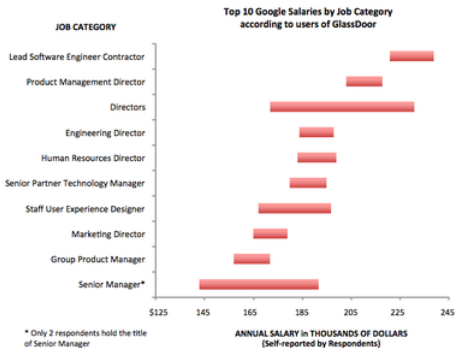


Figure 13: *Top 10 Google salaries by job category: Interval chart.*  
[https://junkcharts.typepad.com/junk\\_charts/2011/10/the-massive-burden-of-pie-charts.html](https://junkcharts.typepad.com/junk_charts/2011/10/the-massive-burden-of-pie-charts.html).



# Bad Graphs: Pie Charts

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

**Bad Graphs**

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

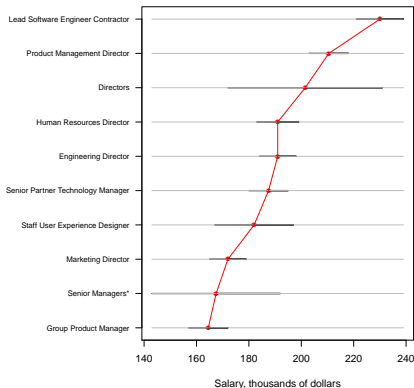


Figure 14: Top 10 Google salaries by job category: Interval chart. Sorted by midpoint of salary range.





# Bad Graphs: Pie Charts

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

**Bad Graphs**

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

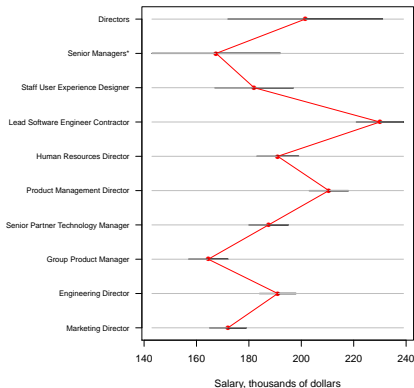


Figure 15: *Top 10 Google salaries by job category: Interval chart. Sorted by salary range.*



# Bad Graphs: Pie Charts

Data Visualization

Dudoit

Motivation

Principles of Data Visualization

Do We Really Need a Graph?

General Considerations

Graphical Perception

**Bad Graphs**

Survey of Data

Visualization Techniques

One Quantitative Variable

Multiple Quantitative Variables

One Qualitative Variable

Multiple Qualitative Variables

Conditional Plots

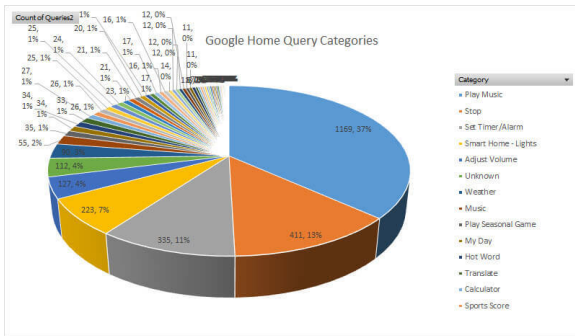


Figure 16: Google Home query categories: Pie chart.

<http://viz.wtf/image/171134950336>. Unreadable. Can't match numbers to categories. What's a better graph?



# Bad Graphs: Pie Charts

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

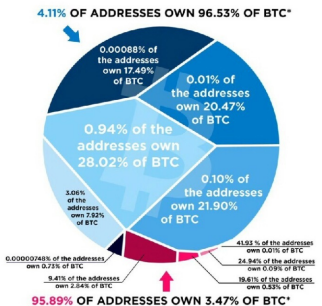
Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

## The Bitcoin Wealth Distribution



\* Data as of September 15th, 2019  
Article and Sources:  
<https://howmuch.net/articles/bitcoin-wealth-distribution>  
<https://bitcoincept.com/>

howmuch.net

Figure 17: Bitcoin wealth distribution: Pie chart.

<http://viz.wtf/image/166329900475>. What's the message?  
How to compare shapes and areas? Without text, pie uninformative.  
What's a better graph?



# Bad Graphs: Pie Charts

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data  
Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

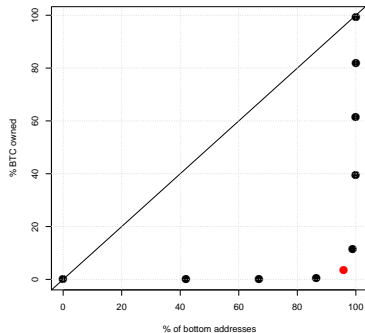
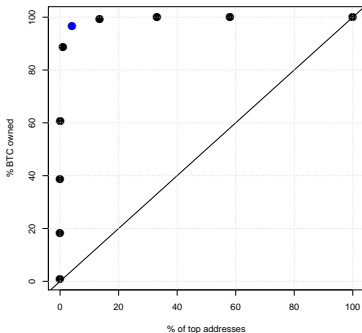


Figure 18: *Bitcoin wealth distribution: Scatterplot.*



# Bad Graphs: Multilevel Donut Charts

Data Visualization

Dudoit

Motivation

Principles of Data Visualization

Do We Really Need a Graph?

General Considerations

Graphical Perception

Bad Graphs

Survey of Data Visualization Techniques

One Quantitative Variable

Multiple Quantitative Variables

One Qualitative Variable

Multiple Qualitative Variables

Conditional Plots

Goldman Sachs open job listings  
Distribution of active job listings as of 9/14/2017

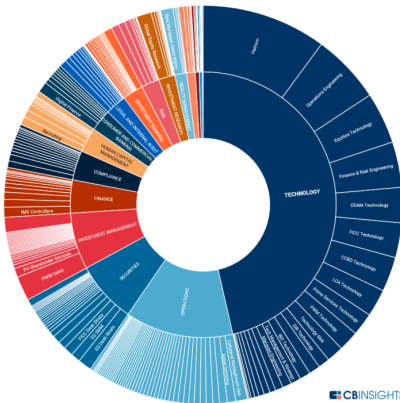


Figure 19: Goldman Sachs job listings: Multilevel donut chart. <https://s3.amazonaws.com/cbi-research-portal-uploads/2017/09/18173935/GStardownjobs>. What's the message? Unreadable. What's a better graph?



# Bad Graphs: Wordclouds

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

**Bad Graphs**

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

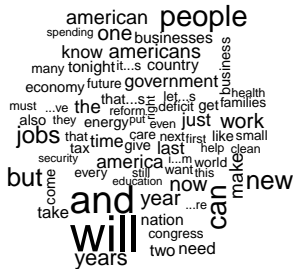


Figure 20: *State of the Union speeches 2010 and 2011: Wordcloud.*  
 Frequency of words with at least 15 occurrences. **What's the message? How to compare frequencies of words? What's a better graph?**



# Bad Graphs: Wordclouds

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data  
Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

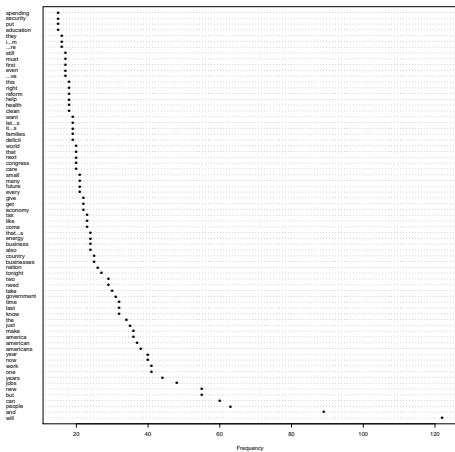


Figure 21: *State of the Union* speeches 2010 and 2011: Dotplot. Frequency of words with at least 15 occurrences.



# Bad Graphs: Wordclouds

Data Visualization

Dudoit

Motivation

Principles of Data Visualization

Do We Really Need a Graph?

General Considerations

Graphical Perception

Bad Graphs

Survey of Data Visualization Techniques

One Quantitative Variable

Multiple Quantitative Variables

One Qualitative Variable

Multiple Qualitative Variables

Conditional Plots



Expressions et hashtags les plus utilisés sur Twitter (Tweets et Retweets) lors des 4 premières journées de mobilisation – Données récoltées avec l'application Talkwalker

Figure 22: *Gilets jaunes: Wordcloud.* Frequency of expressions and hashtags on Twitter for first four days of gilets jaunes movement. How to compare frequencies between days?  
[https://www.lexpress.fr/actualite/societe/gilets-jaunes-ce-qu-en-disent-les-francais\\_2055542.html](https://www.lexpress.fr/actualite/societe/gilets-jaunes-ce-qu-en-disent-les-francais_2055542.html).





# Bad Graphs: Wordclouds

Data Visualization

Dudoit

Motivation

Principles of Data Visualization

Do We Really Need a Graph?

General Considerations

Graphical Perception

**Bad Graphs**

Survey of Data Visualization

Techniques

One Quantitative Variable

Multiple Quantitative Variables

One Qualitative Variable

Multiple Qualitative Variables

Conditional Plots

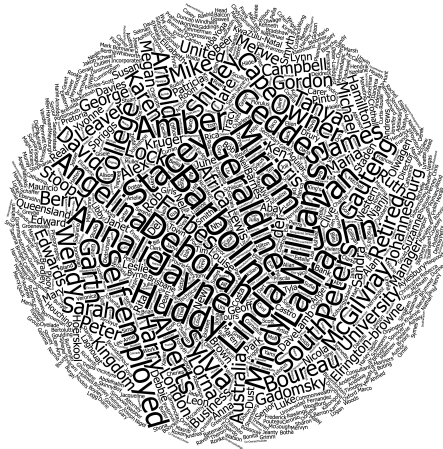


Figure 23: Names: Wordcloud.

[https://www.wordclouds.com/?cloud=names.](https://www.wordclouds.com/?cloud=names)



# Bad Graphs: Wordclouds

Data Visualization

Dudoit

Motivation

Principles of Data Visualization

Do We Really Need a Graph?

General Considerations

Graphical Perception

Bad Graphs

Survey of Data Visualization Techniques

One Quantitative Variable

Multiple Quantitative Variables

One Qualitative Variable

Multiple Qualitative Variables

Conditional Plots



Figure 24: *Business words: Wordcloud.*

<https://www.wordclouds.com/?cloud=business>



# Bad Graphs

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- **Chart junk.** The previous graphs exemplify “chart junk” , i.e., they contain superfluous elements that are not necessary to convey the information contained in the data, but instead distract the viewer from this information or even mask or distort important information.
- **Pie charts.**
  - ▶ Frequency represented by angle/area.
  - ▶ Angles and areas are **hard to perceive and compare**.
  - ▶ Pie charts quickly become **unreadable** for more than a handful of values.
  - ▶ Listing the values is often better – they are actually often added to a pie chart anyway!
  - ▶ How to select **order of categories?**
  - ▶ **Not amenable to comparing distributions**; side-by-side comparisons not effective.



# Bad Graphs

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- ▶ Hard to extend to multiple variables.
- ▶ A lot of **junk** often added to pie charts, e.g., thickness, slice explosion.
- Wordclouds/tag clouds.
  - ▶ Frequency represented by font size.
  - ▶ Neither area nor height corresponds to frequency of words.
  - ▶ How do longer words compare with shorter words?
  - ▶ How are capital letters handled?
  - ▶ How to calculate **relative difference in frequency** between two words?
  - ▶ How are the **words ordered** within the cloud (alphabetical, frequency)?
  - ▶ **Not amenable to comparing distributions**; side-by-side comparisons not effective.
  - ▶ How to **extend** to multiple variables?
  - ▶ A lot of **junk** often added to word clouds.



# Bad Graphs

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Barcharts/barplots. Better.
  - ▶ Based on length and position on common aligned scale.
  - ▶ Add an irrelevant dimension (thickness of bar).
  - ▶ How to select order of categories?
  - ▶ Not readily amenable to comparisons.
  - ▶ Extension to multiple variables problematic.
- Dotcharts/dotplots. (And interval charts.) Even better.
  - ▶ Based on length and position on common aligned scale.
  - ▶ Display only the relevant information.
  - ▶ How to select order of categories?
  - ▶ More amenable to comparisons and extensions to multiple variables.



# Gapminder

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

Gapminder. (<https://www.gapminder.org>)

- We will use data from Gapminder to **reason through the process of data visualization**, e.g., population, population density, life expectancy, income for each country.
- Note that in this case we have a **census**, i.e., there is no sampling involved <sup>3</sup>.
- Gapminder is a Swedish foundation co-created in 2005 by Hans Rosling (Professor of International Health at Karolinska Institute) and family members.
- *“Gapminder is a fact tank, not a think tank.”*  
*“Gapminder measures ignorance about the world.”*  
*“Gapminder makes global data easy to use and understand.”*  
*“Gapminder promotes Factfulness, a new way of thinking.”*



# Gapminder

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Gapminder developed [Trendalyzer](#), a [data visualization software](#) providing [dynamic and interactive graphics](#) of data compiled by organizations such as the United Nations and the World Bank (acquired by Google in 2007).

---

<sup>3</sup>Some of the data could be estimates, but we won't concern ourselves with this at this point.



# Gapminder

Data Visualization

Dudoit

Motivation

Principles of Data Visualization

Do We Really Need a Graph?

General Considerations  
Graphical Perception  
Bad Graphs

Survey of Data Visualization Techniques

One Quantitative Variable  
Multiple Quantitative Variables  
One Qualitative Variable  
Multiple Qualitative Variables  
Conditional Plots

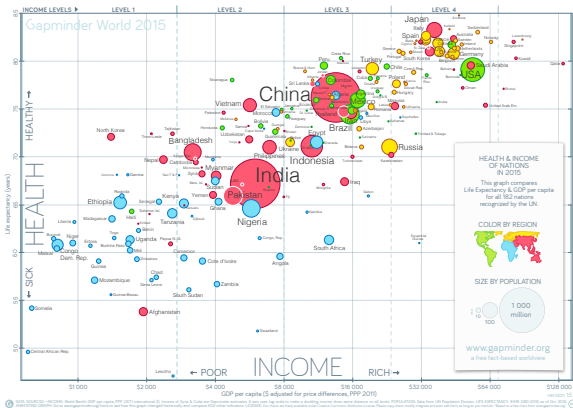


Figure 25: Gapminder: World Poster 2015. "How Does Income Relate to Life Expectancy? Short answer - Rich people live longer." Bubble chart with four variables displayed in 2D.





# Software

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Most of the plots below are produced with Python's seaborn library, using default arguments.
- **Default settings** typically do not correspond to the most basic version of the plot, but rather **impose many decisions** on the plot, e.g., color, legend, ordering. **Experiment with different settings** to make sure you get the plot you want.
- Seaborn tutorial:  
<https://seaborn.pydata.org/tutorial.html>.  
Each function has many arguments to customize the plots. As usual, consult documentation.
- Datasets available at:  
<https://github.com/mwaskom/seaborn-data>.  
E.g. Titanic survival dataset, Fisher's iris dataset.



# One Quantitative Variable

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception  
Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

*How would you visualize life expectancy in 2018 over all countries?*

<b>count</b>	182.000000
<b>mean</b>	72.726374
<b>std</b>	7.237996
<b>min</b>	51.100000
<b>25%</b>	67.150000
<b>50%</b>	74.100000
<b>75%</b>	78.075000
<b>max</b>	84.200000



# Stem-And-Leaf Plots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

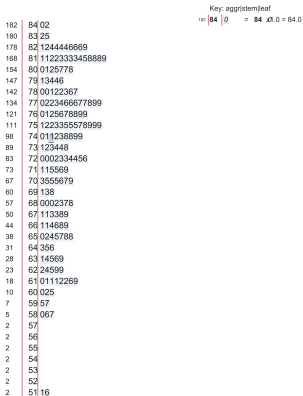


Figure 26: *Life expectancy, 2018.*



# Stripplots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

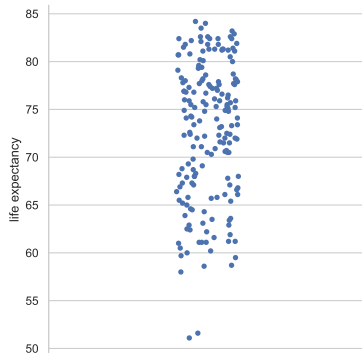
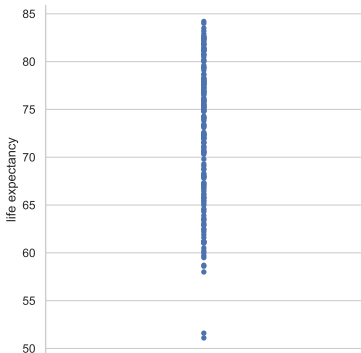


Figure 27: *Life expectancy, 2018*. Right: Jittering, i.e., adding random noise, to avoid overplotting.



# Histograms

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

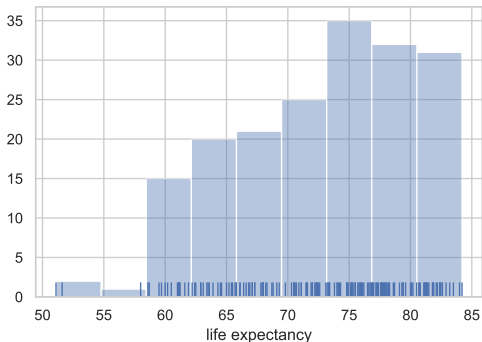


Figure 28: *Life expectancy, 2018.*



# Histograms

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception  
Bad Graphs

Survey of  
Data  
Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

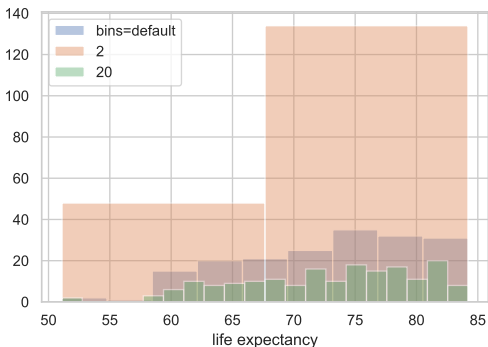


Figure 29: *Life expectancy, 2018*. Different numbers of bins.



# Density Plots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

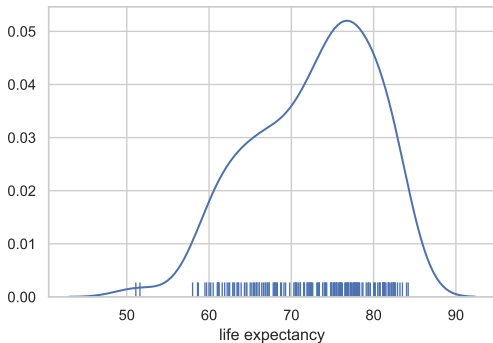


Figure 30: *Life expectancy, 2018.*



# Density Plots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

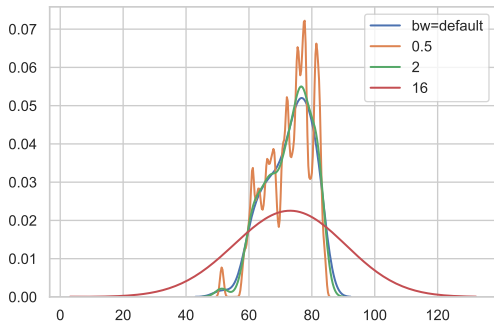


Figure 31: *Life expectancy, 2018*. Different bandwidths.





# Boxplots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

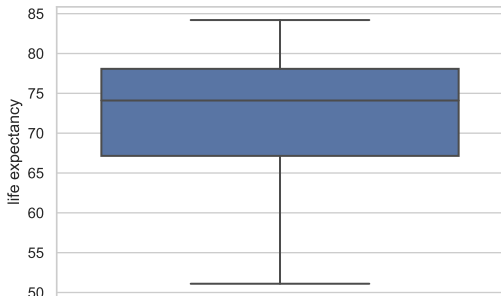


Figure 32: *Life expectancy, 2018.*



# One Quantitative Variable and One Qualitative Variable

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

*How would you visually compare life expectancy between regions?*

```
In [52]: (gm2018.groupby('six_regions'))['life expectancy'].describe()
```

Out[52]:

	count	mean	std	min	25%	50%	75%	max
<b>six_regions</b>								
<b>america</b>	33.0	75.827273	3.774360	64.5	73.400	76.10	78.600	82.2
<b>east_asia_pacific</b>	26.0	72.557692	6.833399	61.1	68.100	71.55	77.275	84.2
<b>europa_central_asia</b>	49.0	77.969388	4.047746	70.5	75.200	78.00	81.500	83.5
<b>middle_east_north_africa</b>	19.0	75.689474	4.642544	67.1	74.050	76.90	78.050	82.4
<b>south_asia</b>	8.0	71.675000	6.652121	58.7	68.825	72.45	75.550	80.1
<b>sub_saharan_africa</b>	47.0	64.157447	4.852599	51.1	61.150	63.90	66.850	74.9



# Stripplots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

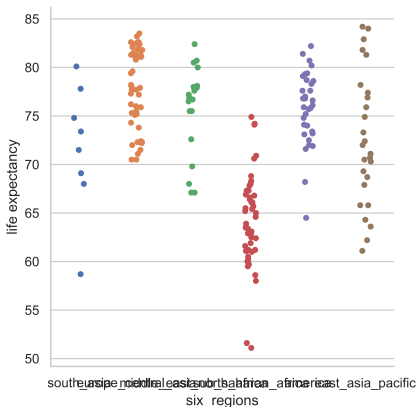


Figure 33: *Life expectancy by region, 2018.*



# Histograms

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

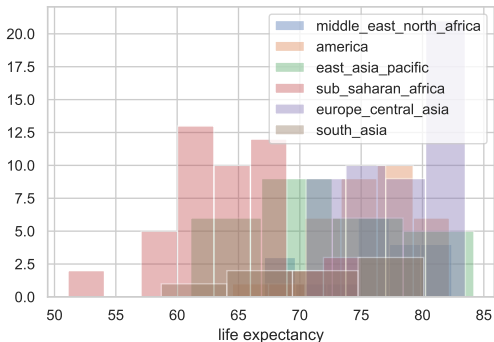


Figure 34: *Life expectancy by region, 2018.*



# Density Plots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

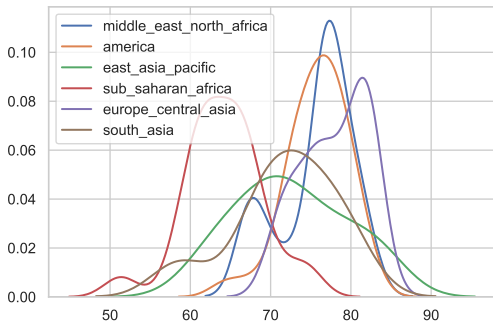


Figure 35: *Life expectancy by region, 2018.*



# Boxplots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

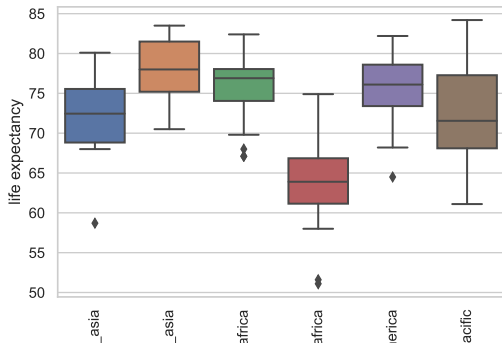


Figure 36: *Life expectancy by region, 2018.*



# Violin Plots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

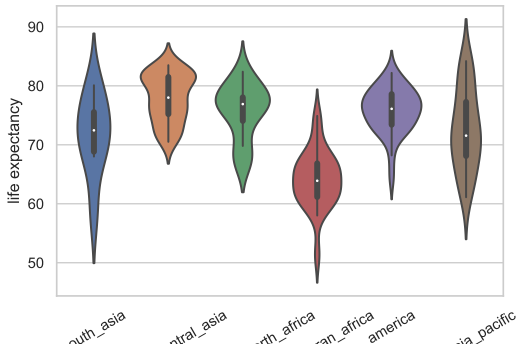


Figure 37: *Life expectancy by region, 2018.*



# Log-Transformation

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

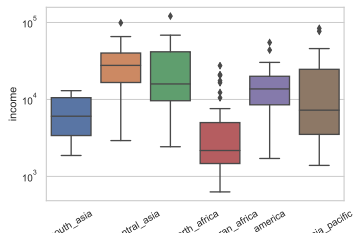
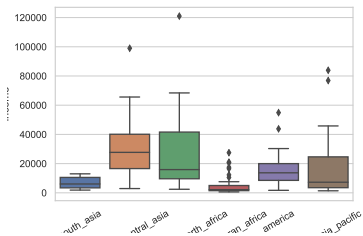


Figure 38: *Income, 2018*. Left: Income (GDP/capita, inflation-adjusted \$). Right: Log-transformed income.





# Time Series

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

**One  
Quantitative  
Variable**

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

*How did life expectancy vary between 1800 and 2018?*



# Time Series

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

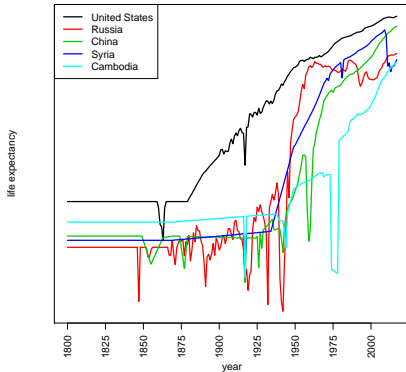


Figure 39: *Life expectancy over time for five countries.*



# Time Series

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

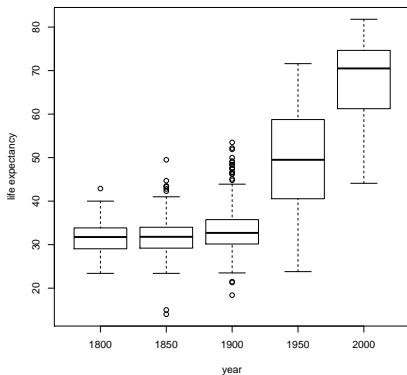


Figure 40: *Life expectancy over time.*



# One Quantitative Variable: Summary

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

Displaying and comparing marginal distributions for quantitative data.

- **Stem-and-leaf plots.**
  - ▶ Simple pen-and-paper method for visualizing the distribution of all of a handful of values.
  - ▶ Not amenable to comparisons between distributions.
  - ▶ No reason to use these days.
- **Stripcharts/Stripplots.** (Sometimes referred to as dotcharts/dotplots, related to rug plots.)
  - ▶ Effective for visualizing the distribution of all of a moderate number of values.
  - ▶ Can use side-by-side stripplots to compare multiple distributions.
- **Histograms.**
  - ▶ Classical method for displaying a single distribution.



# One Quantitative Variable: Summary

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- ▶ Sensitive to bin width and bin boundaries.
- ▶ Cannot easily display and compare multiple distributions.
- **Density plots.**
  - ▶ Based on kernel density estimation (cf. smoothing).
  - ▶ Sensitive to bandwidth, but methods available to select bandwidth.
  - ▶ Effective for displaying and comparing multiple distributions.
- **Boxplots.** (A.k.a., box-and-whiskers plots.)
  - ▶ Summarize distribution by only 5 numbers (+ outliers): Median, upper and lower-quartiles, whiskers at 1.5 times inter-quartile range (IQR) above and below upper and lower-quartiles, respectively.
  - ▶ Possible loss of information, e.g., multimodality.
  - ▶ Effective for displaying and comparing multiple distributions, especially with notches.



# One Quantitative Variable: Summary

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

**One  
Quantitative  
Variable**

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Violin plots.
  - ▶ Trendy hybrids of boxplots and density plots.
  - ▶ Redundant (twice the density plot!), unless plot different densities on each side.
  - ▶ Same limitations and issues as with boxplots and density plots.
  - ▶ Cannot compare densities as readily as with standard density plots.



# Multiple Quantitative Variables

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

**Multiple  
Quantitative  
Variables**

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

*How would you visually examine the relationship between life expectancy and income in 2018 over all countries?*



# Scatterplots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

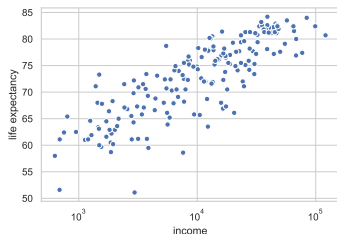
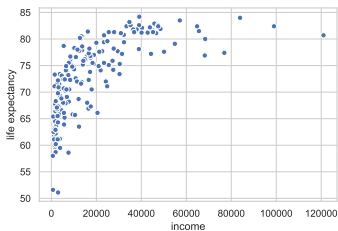


Figure 41: *Life expectancy vs. income, 2018*. Left: Income (GDP/capita, inflation-adjusted \$). Right: Log-transformed income.





# Scatterplots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

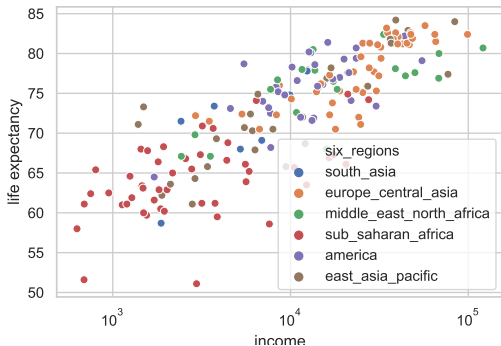


Figure 42: *Life expectancy vs. income, colored by region, 2018.*



# Bubble Charts

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

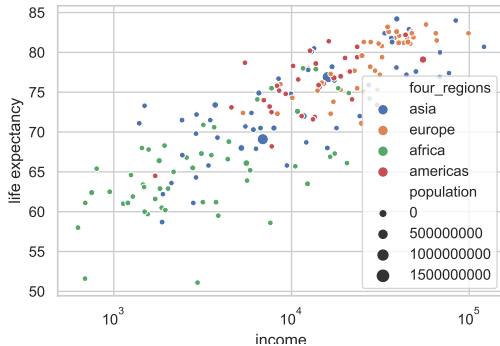


Figure 43: Life expectancy vs. income, colored by region and with area of bubbles representing population, 2018.



# Mean-Difference Plots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

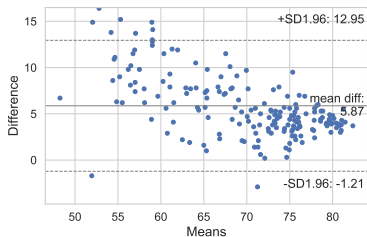
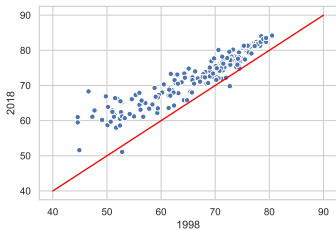


Figure 44: *Life expectancy, 2018 vs. 1998.*



# Scatterplot Matrices

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

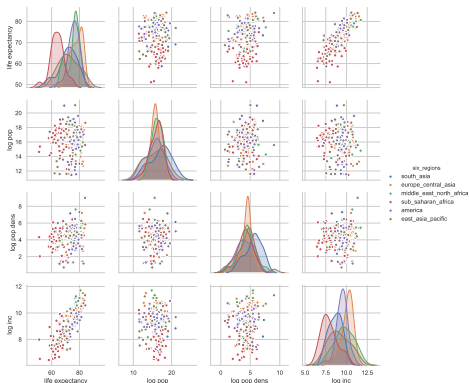


Figure 45: *Life expectancy, population, population density, and income, by region, 2018.*



# Scatterplot Matrices

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

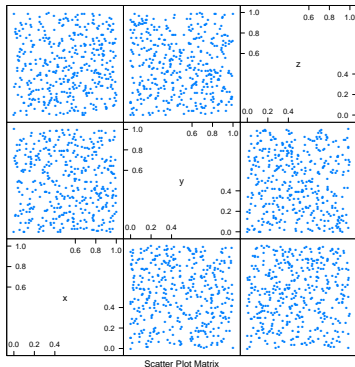


Figure 46: *RANDU* RNG. Triples of successive numbers.



# 3D Scatterplots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

**Multiple  
Quantitative  
Variables**

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

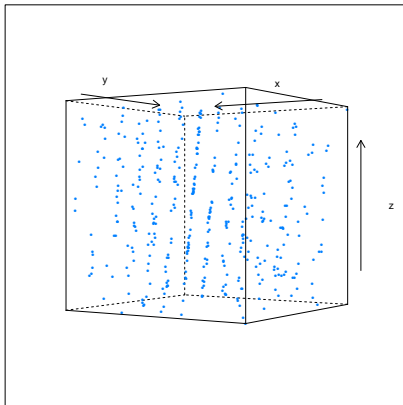


Figure 47: *RANDU* RNG. Triples of successive numbers.



# RANDU RNG

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

RANDU random number generator. (R. Ihaka, <https://www.stat.auckland.ac.nz/~ihaka/120/Lectures/lecture27.pdf>.)

- The dataset consists of 400 triples of successive numbers produced by the RANDU random number generator (RNG).
- The consecutive triples produced by RANDU are constrained to lie on a series of parallel planes which cut through the unit cube.
- The planes are not aligned with the sides of the unit cube and so do not show up in any of the panels of a scatterplot matrix.



# Overplotting

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

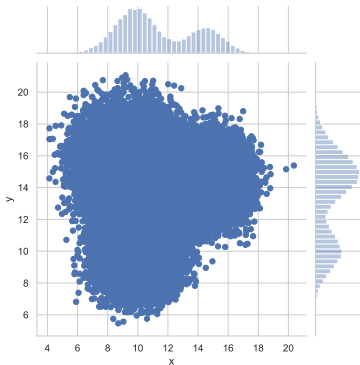


Figure 48: *Simulated data,  $n = 60,000$ : Scatterplot.*





# Overplotting: Hexagonal Binning

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

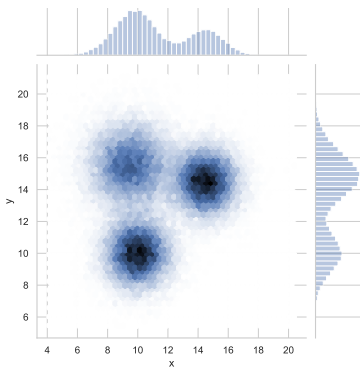


Figure 49: *Simulated data,  $n = 60,000$ : Hexagonal binning.*



# Overplotting: Scatterplot Smoothing

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

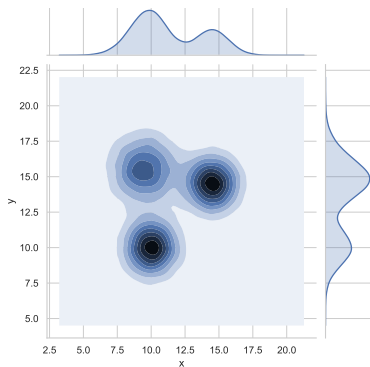


Figure 50: *Simulated data,  $n = 60,000$ : Scatterplot smoothing.*



# Multiple Quantitative Variables: Summary

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

## Displaying joint distributions for quantitative data.

- While density plots and boxplots are useful for comparing two or more **marginal distributions** (e.g., in terms of location and scale), they do not provide any information about **joint distributions** and, in particular, associations between two variables.
- **Scatterplots and scatterplot matrices.**
  - ▶ Useful for examining linear association between two variables.
  - ▶ Can extend beyond two variables by using color and plotting symbol area, as in bubble charts.
  - ▶ However, can miss important higher-dimensional patterns (cf. RANDU example).
- **Mean-difference plots.**



# Multiple Quantitative Variables: Summary

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- ▶ Rotated and scaled version of scatterplot.
- ▶ Better for looking at differences vs. associations.
- **Bubble charts.** A bubble chart is a type of scatterplot that displays one or two extra dimensions using area and color.
- **Parallel coordinates plots.**
  - ▶ Natural for visualizing time series data, i.e., same variable measured across time.  
Cf. Train schedules.
  - ▶ Can also be used for visualizing the relationship between multiple variables, but trickier: Each line corresponds to an observation and each axis to a variable.
  - ▶ Three important considerations, that can affect interpretation of the plot: The order, the rotation, and the scaling of the axes.



# Overplotting

Data  
Visualization

Dudoit

Motivation

Principles of  
Data  
Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data  
Visualization  
Techniques

One  
Quantitative  
Variable

**Multiple  
Quantitative  
Variables**

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

Overplotting issues can be reduced by the following approaches.

- Changing plotting symbol.
- Jittering, i.e., adding random noise.
- Smoothing.
- Hexagonal binning.



# Qualitative Variables

*How would you visualize the 2017 UK election results?*

Number of seats for each of 13 parties.

Party MPs

0 CON 318

1 LAB 261

2 SNP 35

3 LIB DEM 12

4 DUP 10

5 SF 7

6 PC 4

7 GREEN 1

8 IND 1

9 OTHER 1

10 UKIP 0

11 SDLP 0

12 UUP 0

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots



# Pie Charts

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

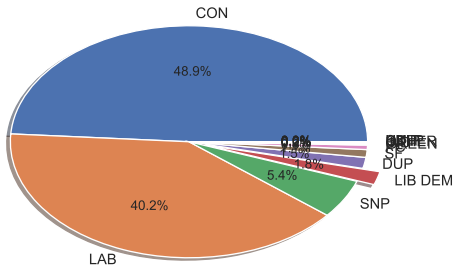


Figure 51: *UK Election Results 2017*. Number of seats for each of 13 parties.



# Barplots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

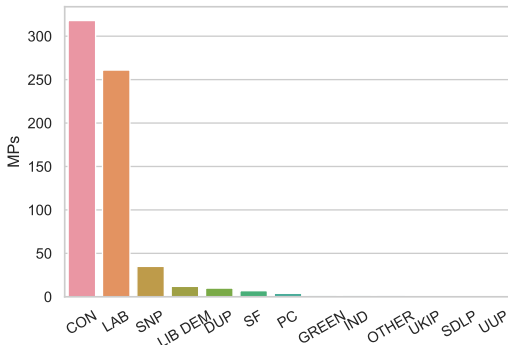


Figure 52: *UK Election Results 2017*. Number of seats for each of 13 parties.





# Dotplots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

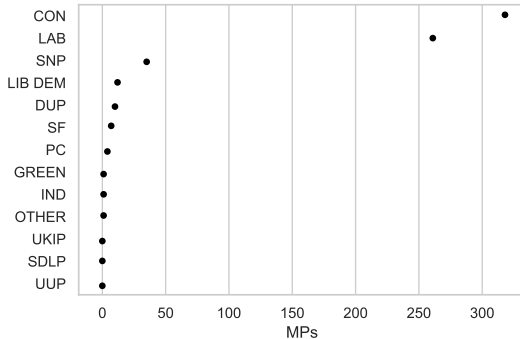


Figure 53: *UK Election Results 2017*. Number of seats for each of 13 parties.



# Lollipop Plots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

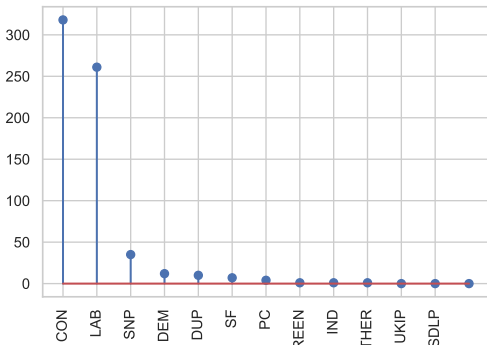


Figure 54: *UK Election Results 2017*. Number of seats for each of 13 parties.



# One Qualitative Variable: Summary

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Pie charts.
  - ▶ Frequency represented by angle/area.
  - ▶ Angles and areas are **hard to perceive and compare**.
  - ▶ Pie charts quickly become **unreadable** for more than a handful of values.
  - ▶ Listing the values is often better – they are actually often added to a pie chart anyway!
  - ▶ How to select **order of categories**?
  - ▶ **Not amenable to comparing distributions**; side-by-side comparisons not effective.
  - ▶ **Hard to extend** to multiple variables.
  - ▶ A lot of **junk** often added to pie charts, e.g., thickness, slice explosion.
- Wordclouds/tag clouds.
  - ▶ Frequency represented by font size.



# One Qualitative Variable: Summary

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- ▶ Neither area nor height corresponds to frequency of words.
- ▶ How do longer words compare with shorter words?
- ▶ How are capital letters handled?
- ▶ How to calculate **relative difference in frequency** between two words?
- ▶ How are the **words ordered** within the cloud (alphabetical, frequency)?
- ▶ **Not amenable to comparing distributions**; side-by-side comparisons not effective.
- ▶ How to **extend** to multiple variables?
- ▶ A lot of **junk** often added to word clouds.
- **Barcharts/barplots.**
  - ▶ Based on **length and position on common aligned scale.**
  - ▶ Add an **irrelevant dimension** (thickness of bar).
  - ▶ How to select **order of categories**?
  - ▶ Not readily amenable to **comparisons.**



# One Qualitative Variable: Summary

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- ▶ Extension to multiple variables problematic.
- Dotcharts/dotplots. (And interval charts.)
  - ▶ Based on length and position on common aligned scale.
  - ▶ Display only the relevant information.
  - ▶ How to select order of categories?
  - ▶ More amenable to comparisons and extensions to multiple variables.
- Lollipop plots.
  - ▶ Similar to dotcharts/dotplots (with added stem) and barcharts/barplots.
  - ▶ Stem is redundant.
  - ▶ How to select order of categories?
  - ▶ Not readily amenable to comparisons.
  - ▶ Extension to multiple variables problematic.



# Multiple Qualitative Variables

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

*How would you display survival data on the Titanic according to class, gender, and age?*

```
In [107]: pd.crosstab(index=titanic['survived'], columns=[titanic['class'], titanic['who']], margins=True)
```

Out[107]:

class	First			Second			Third			All	
	child	man	woman	child	man	woman	child	man	woman		
survived											
0	1	77	2	0	91	6	33	281	58	549	
1	5	42	89	19	8	60	25	38	56	342	
All	6	119	91	19	99	66	58	319	114	891	



# Barplots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

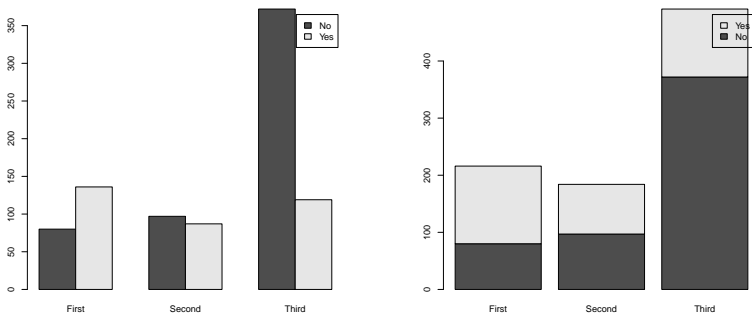


Figure 55: *Titanic: Survival by class.*



# Dotplots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

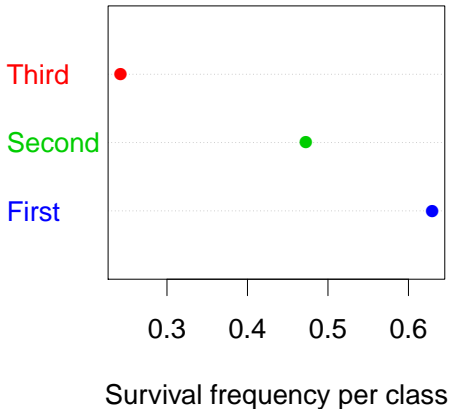


Figure 56: *Titanic: Survival by class.*





# Dotplots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One

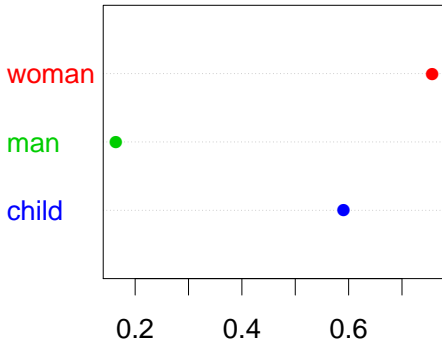
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots



Survival frequency per gender/age

Figure 57: *Titanic: Survival by gender/age.*



# Dotplots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

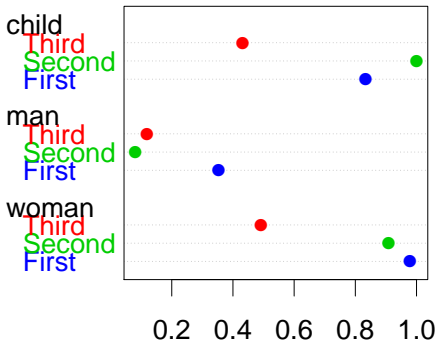
One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots



Survival frequency per class and gender,

Figure 58: *Titanic: Survival by class and gender/age.*



# Mosaic Plots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

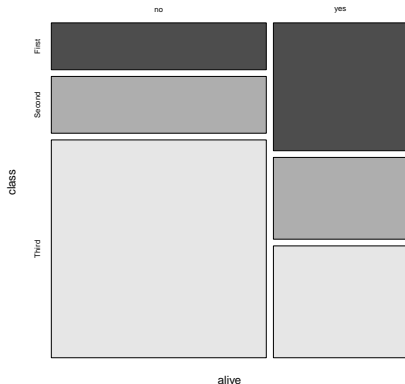


Figure 59: *Titanic: Survival and class.*



# Mosaic Plots

Data Visualization

Dudoit

Motivation

Principles of Data Visualization

Do We Really Need a Graph?

General Considerations

Graphical Perception

Bad Graphs

Survey of Data Visualization Techniques

One Quantitative Variable

Multiple Quantitative Variables

One Qualitative Variable

Multiple Qualitative Variables

Conditional Plots

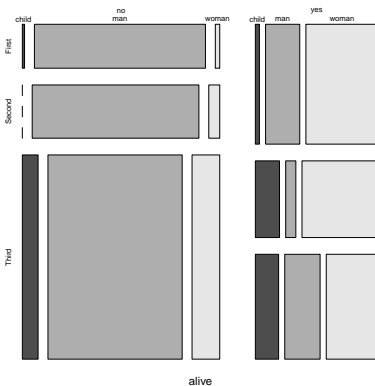


Figure 60: *Titanic: Survival, class, and gender/age.*



# Multiple Qualitative Variables: Summary

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

The following types of plots are used to represent **conditional distributions** for multiple categorical variables or counts for hierarchical categories.

- **Multilevel donut/pie/sunburst plots.**
  - ▶ Same or worse perception issues as with univariate pie charts.
  - ▶ Which variable to choose for “outer” layer?
- **Barcharts/barplots.**
  - ▶ For two categorical variables, a barchart/barplot displays the counts (or percentages) for each category of the second variable within each category of the first variable., i.e., conditional distribution of second variable given first.
  - ▶ Which variable to choose as “first”?
  - ▶ In a **side-by-side barplot**, the frequencies for the second variable are displayed as juxtaposed bars.



# Multiple Qualitative Variables: Summary

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- ▶ In a **stacked/segmented barplot**, the bars for the second variable are staked, so that their total height is the total count for the category of the first variable or 100 percent.
- ▶ Hard to compare frequencies between categories of first variable with both types of barplots.
- ▶ Hard to compare frequencies of second variable within categories of first variable with stacked barplot.
- ▶ Circular barcharts/barplots: Eye-catching, but even harder to compare frequencies.
- **Treemap**. The hierarchical or conditional frequencies are represented using nested figures, usually rectangles.
- **Mosaic plots**.
  - ▶ A mosaic plot is a graphical display of the counts in a **contingency table** (a.k.a., cross-tabulation or crosstab), where each cell is represented by a tile (i.e., rectangle) whose area is proportional to the cell frequency.



# Multiple Qualitative Variables: Summary

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

**Multiple  
Qualitative  
Variables**

Conditional  
Plots

- ▶ Color and shading of the tiles can be used to represent unusually large or small counts, the sign and magnitude of residuals (deviations) for particular models (e.g., independence).
- ▶ For two categorical variables, the width of each tile is proportional to the marginal frequency of the category for the first variable and the height of the tile to the conditional frequency of the category for the second variable given the first.
- ▶ Can be hard to read mosaic plots for more than two variables.



# Conditional Plots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of

Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

## Conditional plots/coplots/faceting/panels/small multiples.

- Collection of plots, where each plot represents the **conditional distribution** of one or more variables given a conditioning variable.
- Each plot corresponds to a value or set of values for the **conditioning variable**. For a quantitative conditioning variable, the ranges are typically chosen so that there are equal numbers of observations in each panel.
- The scales on the axes have to be the same for all panels.
- The colors (and legends) also have to be the same for all panels.
- E.g. Scatterplots of life expectancy vs. income for each of the six world regions.





# Conditional Plots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

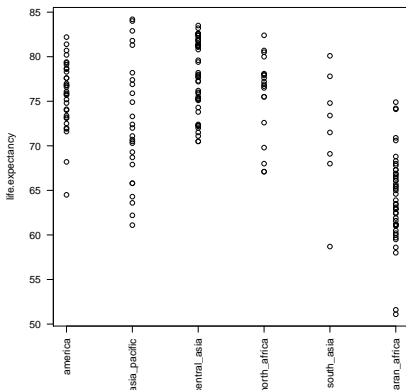


Figure 61: *Life expectancy by region, 2018.*



# Conditional Plots

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization  
Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

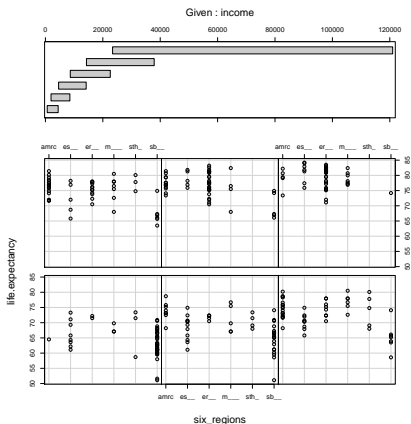


Figure 62: *Life expectancy by region conditioning on income, 2018.*



# References

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization

Do We Really  
Need a Graph?

General  
Considerations

Graphical  
Perception

Bad Graphs

Survey of  
Data

Visualization

Techniques

One  
Quantitative  
Variable

Multiple  
Quantitative  
Variables

One Qualitative  
Variable

Multiple  
Qualitative  
Variables

Conditional  
Plots

- Peter Aldhous. Data visualization: basic principles. <http://paldhous.github.io/ucb/2016/dataviz/week2.html>.
- Ross Ihaka. Statistics 120 – Information Visualisation. <https://www.stat.auckland.ac.nz/~ihaka/120/>.
- Duncan Temple Lang. Data Visualization Workshops. <http://dsi.ucdavis.edu/tag/data-visualization.html>.

- W. S. Cleveland and R. McGill. Graphical perception and graphical methods for analyzing scientific data. *Science*, 229(4716):828–833, 1985.
- A. Gelman, C. Pasarica, and R. Dodhia. Lets practice what we preach: Turning tables into graphs. *The American Statistician*, 56(2):121–130, 2002.
- E. J. Marey. *La Mthode Graphique*. Librairie de l'Académie de Médecine, 1885.



# References

Data  
Visualization

Dudoit

Motivation

Principles of  
Data

Visualization  
Do We Really  
Need a Graph?

General  
Considerations  
Graphical  
Perception  
Bad Graphs

Survey of  
Data  
Visualization  
Techniques

One  
Quantitative  
Variable  
Multiple  
Quantitative  
Variables  
One Qualitative  
Variable  
Multiple  
Qualitative  
Variables  
Conditional  
Plots

- S. S. Stevens. On the psychophysical law. *Psychological Review*, 64(3): 153–181, 1957.
- E. R. Tufte. *The Visual Display of Quantitative Information*. Graphics Press, 2nd edition, 2001.