



BITS Pilani

Pilani | Dubai | Goa | Hyderabad

INTRODUCTION TO DATA SCIENCE

MODULE # 7 : DATA VISUALIZATION

IDS Course Team

BITS Pilani

The instructor is gratefully acknowledging
the authors who made their course
materials freely available online.

TABLE OF CONTENTS

1 DATA VISUALIZATION

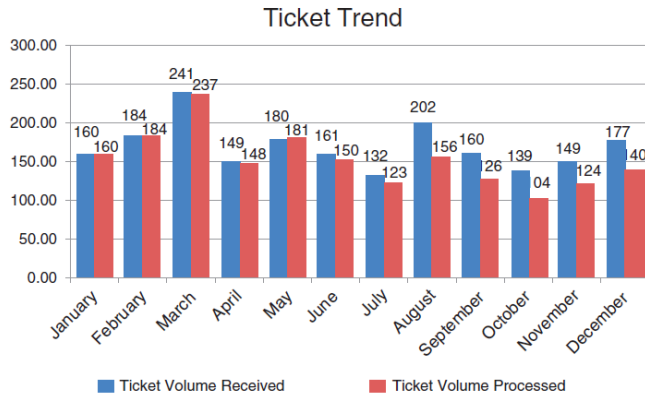
2 VISUALIZATION CONTEXT

3 STORYTELLING STRATEGIES

TYPES OF DATA REPRESENTATIONS

- Presentation
 - ▶ Uses data visuals to communicate
 - ▶ Two roles : Presenter and Audience
 - ▶ Communicate and Persuade
- Visualization
 - ▶ Use visuals to think.
 - ▶ Involves people trying to answer questions.

CASE 1



CASE 1



Better Visualization

Please approve the hire of 2 FTEs

to backfill those who quit in the past year

Ticket volume over time

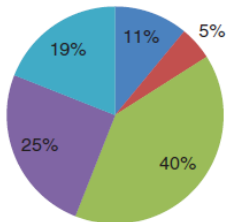


Data source: XYZ Dashboard, as of 12/31/2014. | A detailed analysis on tickets processed per person and time to resolve issues was undertaken to inform this request and can be provided if needed.

Survey Results

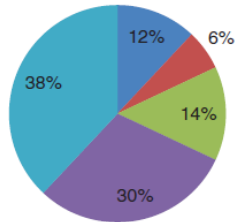
PRE: How do you feel about doing science?

■ Bored ■ Not great ■ OK ■ Kind of interested ■ Excited



POST: How do you feel about doing science?

■ Bored ■ Not great ■ OK ■ Kind of interested ■ Excited



CASE 2

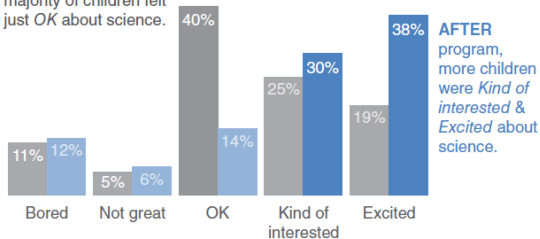


Better Visualization

Pilot program was a success

How do you feel about science?

BEFORE program, the majority of children felt just *OK* about science.

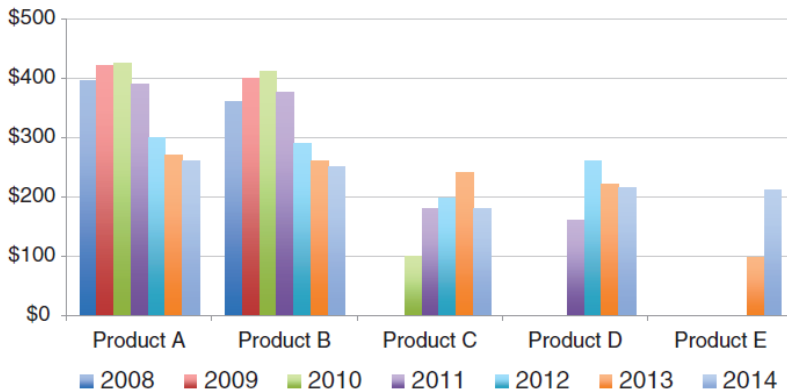


Based on survey of 100 students conducted before and after pilot program (100% response rate on both surveys).

CASE 3



Average Retail Product Price per Year



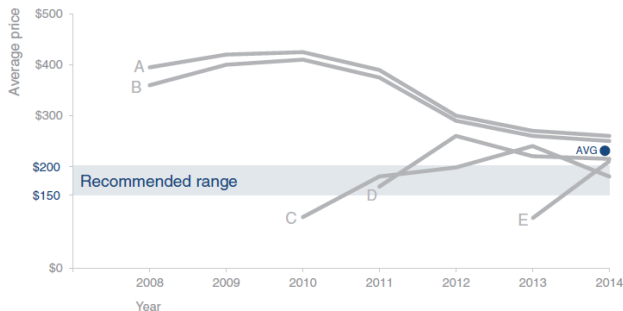
CASE 3



Better Visualization

To be competitive, we recommend introducing our product *below the \$223 average price point* in the **\$150–\$200 range**

Retail price over time by product



DATA VISUALIZATION

- Tell stories with numbers.
- Communicate something to someone using data.

BEN SHNEIDERMAN

“The purpose of visualization is insight, not just picture.”

NEED FOR DATA VISUALIZATION

- What?
 - ▶ Graphical / Visual representation of data.
- Why?
 - ▶ Way to identify patterns, trends and outliers in data.
 - ▶ Helps in making data-driven decisions.

STEPHEN FEW

“Data visualization is the use of visual representations to explore, make sense of, and communicate data.”

DATA VISUALIZATION GOALS

We visualize quantitative data to perform three fundamental tasks in an effort to achieve three essential goals:

Task	Description	Goal
Exploration	Searching for significant facts	Discovery
Sense-making	Examining and making sense of data	Understanding
Communication	Conveying information to the audience	Informed decisions

WHY DO WE REALLY VISUALIZE DATA?

- Because humans can perceive particular quantitative features and perform particular quantitative tasks most effectively when the data is expressed graphically.
- Visual data processing provides optimal support for the following:
 - ▶ Seeing the big picture
 - ▶ Easily and rapidly comparing values
 - ▶ Seeing patterns among values
 - ▶ Comparing patterns

TYPES OF DATA VISUALIZATION

- Based on Interactivity

- ▶ Static Visualizations

- ★ Used in presentations, documents, info-graphics.
 - ★ Requires careful design as it is meant for offline viewing.
 - ★ User cannot adjust the views.
 - ★ Usually focused on a specific data story, users can't go beyond a single view to explore additional stories beyond what's in front of them.
 - ★ The story is specifically captured in an engaging single page layout.

- ▶ Dynamic Visualizations

- ★ Used for exploratory data analysis.
 - ★ Meant for live or online interactions.
 - ★ Users have more viewing options.
 - ★ User gains more control over the display.
 - ★ Users can select specific data points to build a visualized story of their choosing.
 - ★ These visualizations allow the user to be part of the data visualization process by building a story of their choosing.

TYPES OF DATA VISUALIZATION

- Based on Application

- ▶ Common/ General Types of Data Visualization

- * Charts
 - * Maps
 - * Tables
 - * Info-graphics
 - * Graphs
 - * Dashboards
 - * Plots

- ▶ Specific methods

- * Area chart
 - * Box and Whisker plots
 - * Bubble cloud
 - * Cartogram
 - * Circle view
 - * Dot distribution map
 - * Gantt Chart
 - * Heat Map
 - * Highlight table
 - * Histogram
 - * Matrix
 - * Network
 - * Polar Area
 - * Radial Tree
 - * Scatter plot
 - * Stream graph
 - * Text tables
 - * Time-line
 - * Tree map
 - * Wedge Stack graph
 - * Word cloud



USAGE OF DATA VISUALIZATION

- Exploratory Analysis

- ▶ Purpose is to understand the data and figure out what might be noteworthy or interesting to highlight to others.
- ▶ Explore or dig through the data to become familiar with data.
- ▶ Find trends and relationships w.r.t. specific goals.
- ▶ Helps to determine the analyses to apply to data.
- ▶ Eg: hunting for pearls in oysters.

- Explanatory Analysis

- ▶ Explain outcomes or results of analysis. It's about communicating our findings.
- ▶ Tell a story with data.
- ▶ We need to be in the explanatory space when we want to communicate the results of our analysis to our audience.
- ▶ At this stage, we should have a specific story to tell or a specific thing we want to explain.
- ▶ Eg: probably about those two specific pearls.

TABLE OF CONTENTS

- 1 DATA VISUALIZATION
- 2 VISUALIZATION CONTEXT
- 3 STORYTELLING STRATEGIES

VISUALIZATION CONTEXT

COLE KNAFLIC

Success in data visualization does not start with data visualization itself.
... understanding the context sets solid foundation for data visualization creation.

- Understand the context for the need to communicate.
- *W²H* – Who, What, How.

VISUALIZATION CONTEXT

- Who

- ▶ To whom are you communicating?
- ▶ Who is your audience?
 - ★ Knowing their place puts you in a better position for communication.
 - ★ Be specific while identifying audience.
 - ★ Different content for different sets of audience.
 - ★ Avoid general audiences, such as "internal and external stakeholders" or "anyone who might be interested".

- What

- ▶ What do you want your audience to know or do?

- How

- ▶ How can you use data to help make your point?
- ▶ What data is available?
- ▶ Supporting evidence of the story.

AUDIENCE TO KNOW OR ACT

● What

- ▶ Action – What do you want your audience to know or do?
 - ★ Make sure audience care about what you say.
 - ★ You are subject matter expert – unique position to interpret the data and help lead people to understanding and take action.
 - ★ If no action recommendation possible / feasible, then encourage discussion towards one.
- ▶ Mechanism – How you will communicate to your audience?
 - ★ Determines level of control and level of detail.
 - ★ Communication mechanism can be placed on a continuum of live presentation on one end to a written document/email on the other.
 - ★ Three levels of communication mechanisms - Live presentation or Written document or Slideument
- ▶ Tone – How you want your audience to act?
 - ★ Overall tone that you want to set for your communication.
 - ★ Celebrating success? Is topic light-hearted or serious?
 - ★ Lighting a fire to drive action?

COMMUNICATION MECHANISM

LIVE PRESENTATION WRITTEN DOC or EMAIL

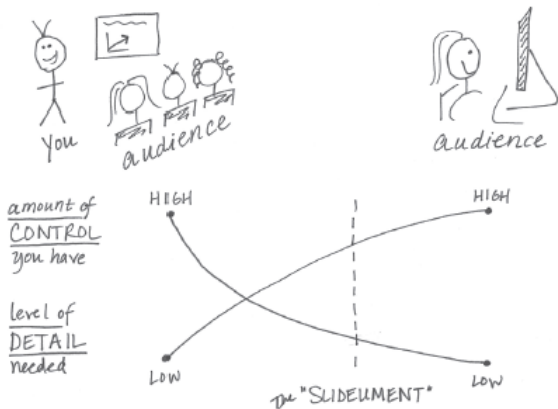


FIGURE 1.1 Communication mechanism continuum

TABLE OF CONTENTS

- 1 DATA VISUALIZATION
- 2 VISUALIZATION CONTEXT
- 3 STORYTELLING STRATEGIES

STORYTELLING STRATEGIES

- The 3-minute story
- Big Idea
- Storyboarding

THE 3-MINUTE STORY

- Within 3 minutes you have to tell audience, what they need to know .
- This is a great way to ensure we are clear on and can articulate the story we want to tell.
- Removes dependence on supporting material like visuals , presentations etc.
- Helps in a situation where we have to give elevator speech, or if our half hour agenda is cut short to ten minutes, or to five.
- Need to know exactly what data is saying.
- Being concise is more challenging than being verbose.

- Boils down to a single most important sentence.
- As per Nancy Duarte (in her book Resonate, 2010), big Idea has three components:
 - ▶ It must articulate your unique point of view.
 - ▶ It must convey what's at stake.
 - ▶ It must be complete sentence.

STORY BOARDING

- Storyboarding helps ensure the communication we craft is on point.
- A storyboard establishes a structure for our communication.
- It is a visual outline of the content we plan to create .
- It can be subject to change as we work through the details.
- Avoid using presentation computer-aided tools in beginning.
- Use whiteboard, post-it notes, paper.
- When possible, get acceptance from your client or stakeholder at this step. It will help ensure that what you're planning is in line with the need.

STORY BOARDING EXAMPLE

Issue:

Kids have bad attitudes about science

Demonstrate Issue:
show student assignment grades over course of year

Ideas for overcoming issue, including pilot program

Describe pilot program - goals, etc.

Show before & after survey data to demonstrate success of program

RECOMMENDATION:
pilot was a success let's expand it we need \$\$\$

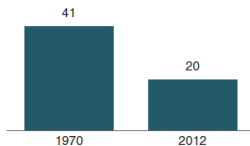
TOOLS FOR DATA VISUALIZATION

- ❶ Simple Text
- ❷ Table
- ❸ Heatmap
- ❹ Scatterplot
- ❺ Line
- ❻ Slopegraph
- ❼ Bar graphs
- ❽ Stacked bar graphs
- ❾ Waterfall
- ❿ Square area

- Use a number or two directly.
- Use the number solely to make it prominent.
- Add a few supporting words to improve clarity.

Children with a "Traditional" Stay-at- Home Mother

*% of children with a married
stay-at-home mother with a
working husband*



Note: Based on children younger than 18. Their mothers are categorized based on employment status in 1970 and 2012.

Source: Pew Research Center analysis of March Current Population Surveys Integrated Public Use Microdata Series (IPUMS-CPS), 1971 and 2013

Adapted from PEW RESEARCH CENTER

20%

of children had a
traditional stay-at-home mom
in 2012, compared to 41% in 1970

FIGURE 2.3 Stay-at-home moms simple text makeover

TABLE

- When we have more numbers to show.
- Interact with our verbal system.
- Used to compare values, communicate to a mixed audience, or multiple different units.
- Design of the table should be at minimum.
- No heavy borders or shading.
- White space or light/gray borders as separator.
- Data should be highlighted.

TABLE



Heavy borders

Group	Metric A	Metric B	Metric C
Group 1	\$X.X	Y%	Z,ZZZ
Group 2	\$X.X	Y%	Z,ZZZ
Group 3	\$X.X	Y%	Z,ZZZ
Group 4	\$X.X	Y%	Z,ZZZ
Group 5	\$X.X	Y%	Z,ZZZ

Light borders

Group	Metric A	Metric B	Metric C
Group 1	\$X.X	Y%	Z,ZZZ
Group 2	\$X.X	Y%	Z,ZZZ
Group 3	\$X.X	Y%	Z,ZZZ
Group 4	\$X.X	Y%	Z,ZZZ
Group 5	\$X.X	Y%	Z,ZZZ

Minimal borders

Group	Metric A	Metric B	Metric C
Group 1	\$X.X	Y%	Z,ZZZ
Group 2	\$X.X	Y%	Z,ZZZ
Group 3	\$X.X	Y%	Z,ZZZ
Group 4	\$X.X	Y%	Z,ZZZ
Group 5	\$X.X	Y%	Z,ZZZ

FIGURE 2.4 Table borders

- Using visual cues in a heatmap.
- A heatmap is a way to visualize data in tabular format, where in place of the numbers, you leverage colored cells that convey the relative magnitude of the numbers.
- Use color saturation to provide visual cues to quickly target the potential points of interest.
- Always include a legend as a subtitle on the heatmap with color corresponding to the conditional formatting color.

HEATMAP



Table

	A	B	C
Category 1	15%	22%	42%
Category 2	40%	36%	20%
Category 3	35%	17%	34%
Category 4	30%	29%	26%
Category 5	55%	30%	58%
Category 6	11%	25%	49%

Heatmap

LOW-HIGH

	A	B	C
Category 1	15%	22%	42%
Category 2	40%	36%	20%
Category 3	35%	17%	34%
Category 4	30%	29%	26%
Category 5	55%	30%	58%
Category 6	11%	25%	49%

FIGURE 2.5 Two views of the same data

SCATTERPLOT

- Show relationship between data.
- Better for categorical data.

SCATTERPLOT



Cost per mile by miles driven

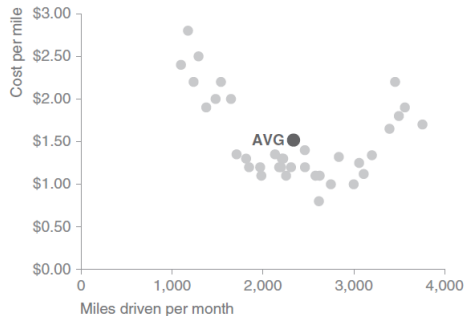


FIGURE 2.6 Scatterplot

Cost per mile by miles driven

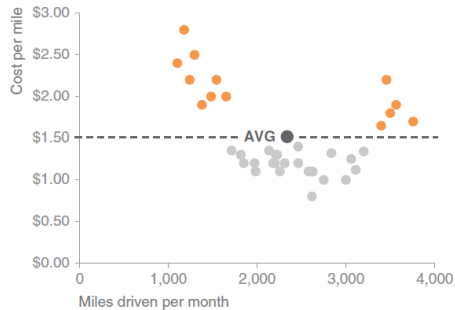
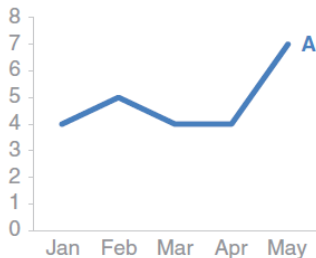


FIGURE 2.7 Modified scatterplot

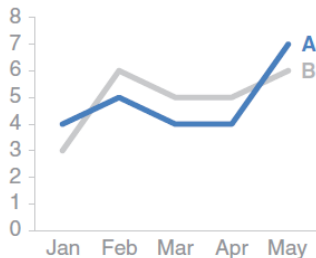
- When we have more numbers to show.
- Interact with our visual system.
- Faster at processing information.
- 4 categories
 - ▶ Points
 - ▶ Lines
 - ▶ Bars
 - ▶ Area

- Plot continuous data.
- The line graph can show a single series of data, two series of data, or multiple series.
- Be consistent in the time points you plot.

Single series



Two series



Multiple series

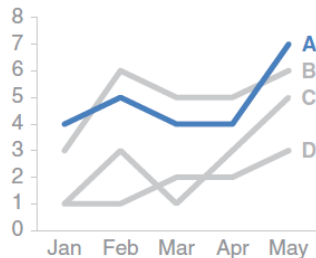


FIGURE 2.8 Line graphs

Passport control wait time

Past 13 months

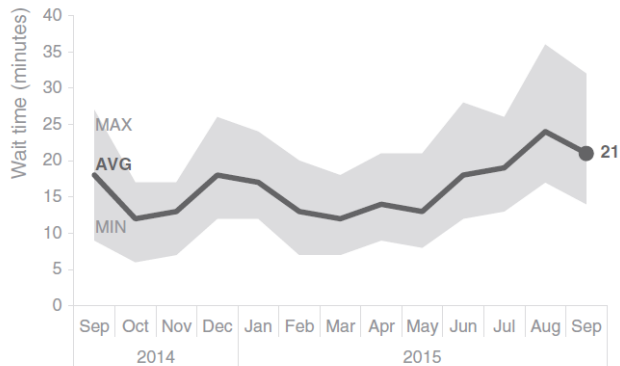


FIGURE 2.9 Showing average within a range in a line graph

- Slopegraphs can be useful when you have two time periods or points of comparison and want to quickly show relative increases and decreases or differences across various categories between the two data points.
- Less overlapping lines.
- Emphasize a single series for highlighting.

Employee feedback over time

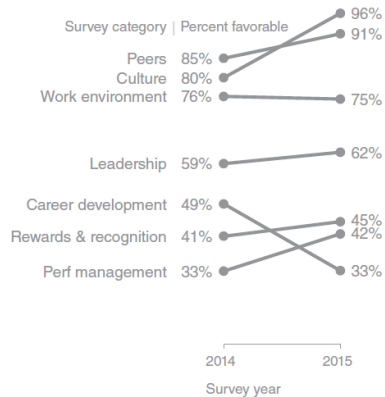


FIGURE 2.10 Slopegraph

Employee feedback over time

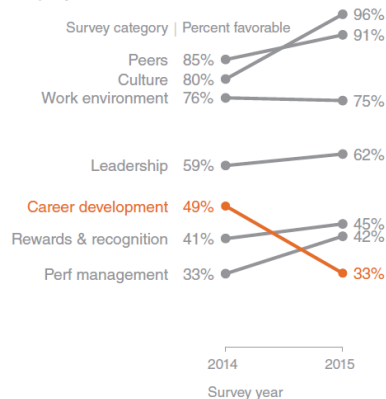


FIGURE 2.11 Modified slopegraph

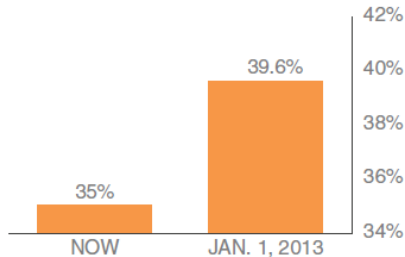
- Compare the end points of the bars, so it is easy to see quickly which category is the biggest, which is the smallest, and also the incremental difference between categories.
- Bar charts should always have a zero baseline.
- Less learning curve for the audience.
- Categorical data or data is organised into groups.

BAR CHARTS



Non-zero baseline: as originally graphed

IF BUSH TAX CUTS EXPIRE
TOP TAX RATE



Zero baseline: as it should be graphed

IF BUSH TAX CUTS EXPIRE
TOP TAX RATE

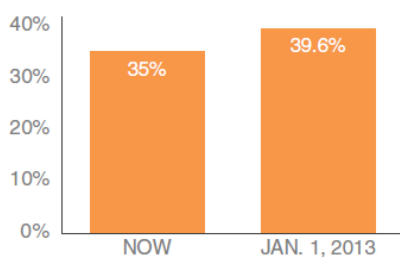
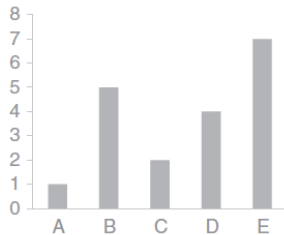


FIGURE 2.13 Bar charts must have a zero baseline

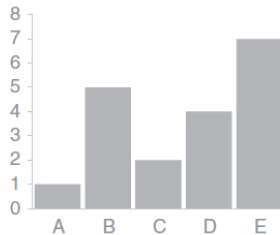
BAR CHARTS



Too thin



Too thick



Just right

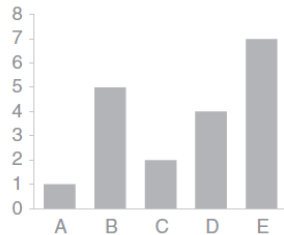


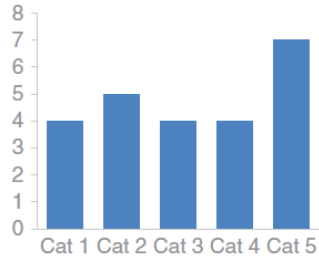
FIGURE 2.14 Bar width

VERTICAL BAR CHARTS

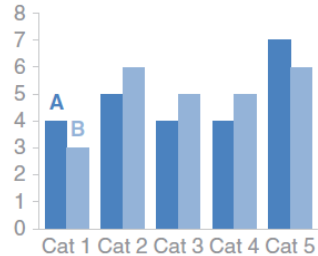
- Ensure that there is visual grouping when using multiple series of vertical bars.
- The relative order of the categorization is important.

VERTICAL BAR CHARTS

Single series



Two series



Multiple series

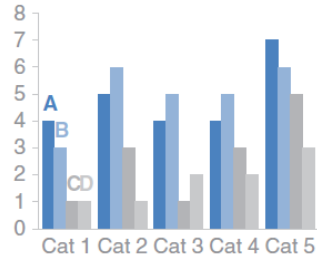


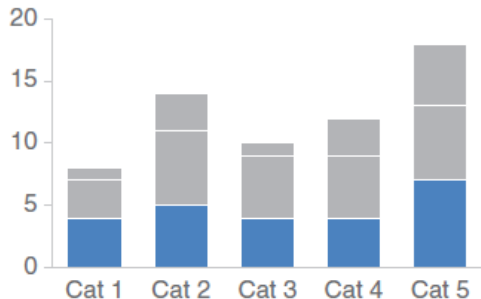
FIGURE 2.15 Bar charts

STACKED VERTICAL BAR CHARTS

- Compare totals across categories and also compare the subcomponent pieces within a given category.
- Visually overwhelming
- Less use of the varied default color schemes.
- The stacked vertical bar chart can be absolute numbers or with each column summing to 100%.
- With 100% stacked bar, include the absolute numbers for each category total for better interpretation of the data.

STACKED VERTICAL BAR CHARTS

Comparing **these** is easy



Comparing **these** is hard

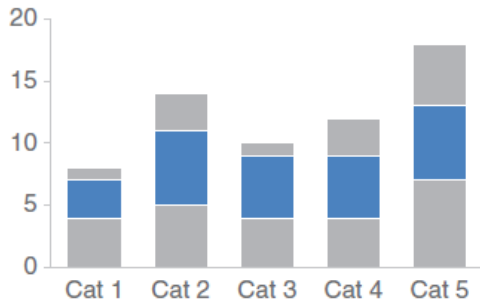


FIGURE 2.16 Comparing series with stacked bar charts

WATERFALL CHARTS

- The waterfall chart can be used to pull apart the pieces of a stacked bar chart to focus on one at a time.
- Show a starting point, incremental additions and deductions, and the resulting ending point.

WATERFALL CHARTS



2014 Headcount math

Though more employees transferred out of the team than transferred in, aggressive hiring means overall headcount (HC) increased 16% over the course of the year.

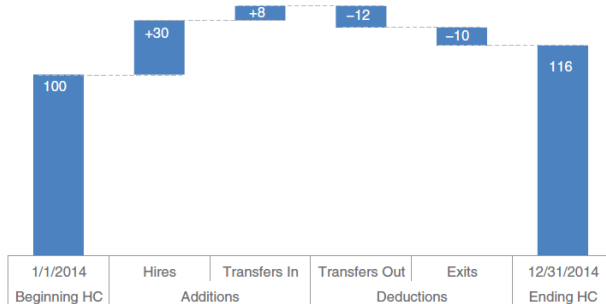


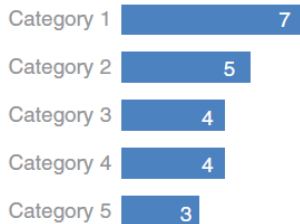
FIGURE 2.17 Waterfall chart

HORIZONTAL BAR CHARTS

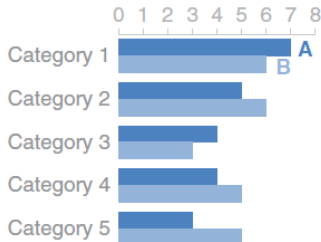
- Categorical data
- Keep categories such that the natural ordering is preserved.
- If the biggest category is the most important, put that first and ordering the rest of the categories in decreasing numerical order.
- If the smallest is most important, put that at the top and order by ascending data values.

HORIZONTAL BAR CHARTS

Single series



Two series



Multiple series

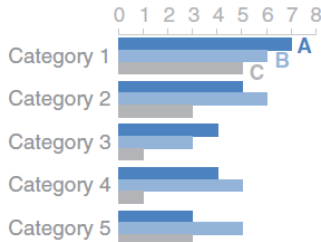


FIGURE 2.18 Horizontal bar charts

STACKED HORIZONTAL BAR CHARTS

- Compare totals across categories and also compare the subcomponent pieces within a given category.
- The stacked vertical bar chart can be absolute numbers or with each row summing to 100%.
- A consistent baseline is available on both the far left and the far right.
- Use to show a scale from negative to positive.

STACKED HORIZONTAL BAR CHARTS

Survey results

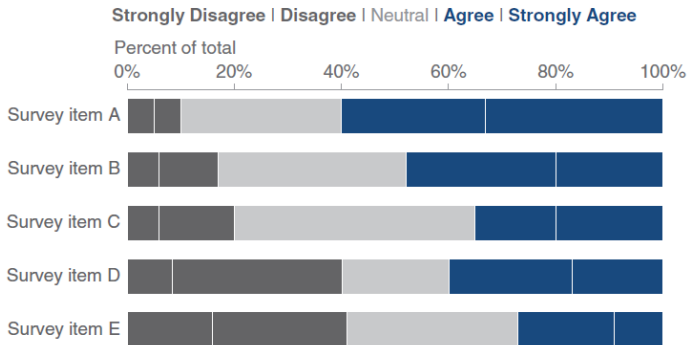


FIGURE 2.19 100% stacked horizontal bar chart

- 2-dimensional

Interview breakdown



FIGURE 2.20 Square area graph

AVOID PIE CHARTS

- Pie Charts to be avoided.
- Use bar charts instead.

Supplier Market Share

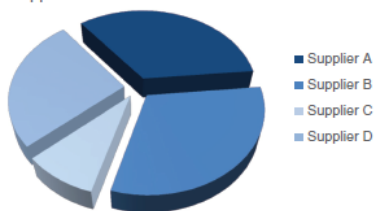


FIGURE 2.21 Pie chart

Supplier Market Share

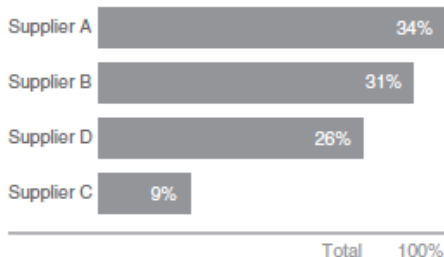


FIGURE 2.23 An alternative to the pie chart

AVOID SECONDARY Y-AXIS

- Label directly.
- Pull apart vertically.

Secondary y-axis

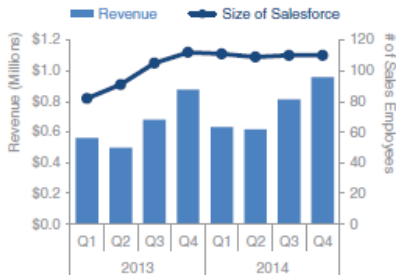


FIGURE 2.26 Secondary y-axis

Alternative 1: label directly



Alternative 2: pull apart vertically

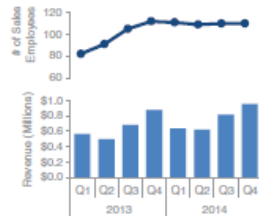


FIGURE 2.27 Strategies for avoiding a secondary y-axis

BETTER GRAPHS

- Reduce clutter
- Alignment of text should be made in context of the other elements on the page.
- Upper-left-most justify the text like title, axis titles, legend.
- Diagonal components should be avoided.
- Include lot of white space.
- Margins should remain free of text and visuals.
- The lack of clear contrast can be a visual clutter.

- Remove chart border.
- Remove gridlines.
- Remove data markers.
- Clean up axis labels.
- Label data directly.
- Leverage consistent color.

PRE-ATTENTIVE ATTRIBUTES

- Employ pre-attentive attributes to create visual hierarchy in the communication.
- **Bold**
- *Italics*
- Color
- Size
- Separate spatially
- Outline
- Underline

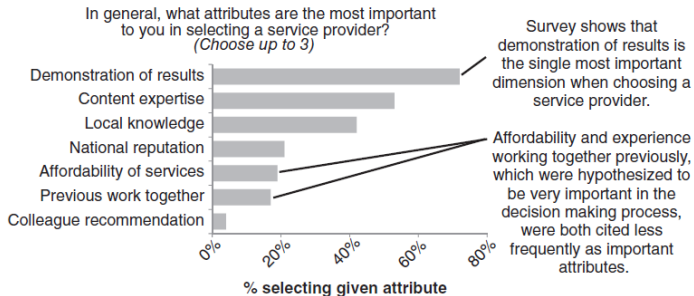
PRE-ATTENTIVE ATTRIBUTES

- Enable the audience to see what we want them to see before they even know they're seeing it!
- A bright blue will typically draw attention more than a muted blue. Both will draw more attention than a light gray.

PRE-ATTENTIVE ATTRIBUTES IN GRAPHS

- Leverage color to draw attention.
- Create a visual hierarchy of information.
- Data labels used sparingly help draw attention.
- Relative size denotes relative importance.
- Use color sparingly.
- Use color consistently.
- Design with colorblind in mind.
- Be thoughtful of tone that color conveys, as color evokes emotion.
- Position on page. Users read in Z pattern.

Demonstrating effectiveness is most important consideration when selecting a provider



Data source: xyz; includes N number of survey respondents. Note that respondents were able to choose up to 3 options.

FIGURE 3.13 Summary of survey feedback

THIS ONE?



Demonstrating effectiveness is most important consideration when selecting a provider

In general, **what attributes are the most important** to you in selecting a service provider?

(Choose up to 3)



Survey shows that **demonstration of results** is the single most important dimension when choosing a service provider.

Affordability and **experience working together previously**, which were hypothesized to be very important in the decision making process, were both cited less frequently as important attributes.

Data source: xyz; includes N number of survey respondents.
Note that respondents were able to choose up to 3 options.

FIGURE 3.14 Revamped summary of survey feedback

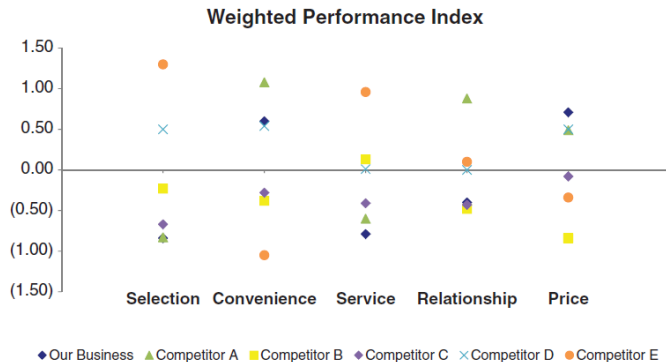


FIGURE 3.15 Original graph

THIS ONE?



Performance overview

■ Our business

- Competitor A
- Competitor B
- Competitor C
- Competitor D
- Competitor E

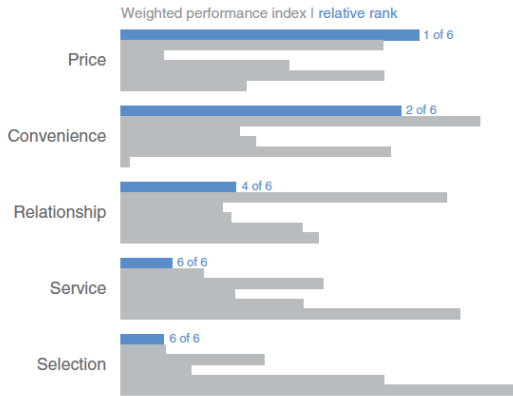


FIGURE 3.16 Revamped graph, using contrast strategically

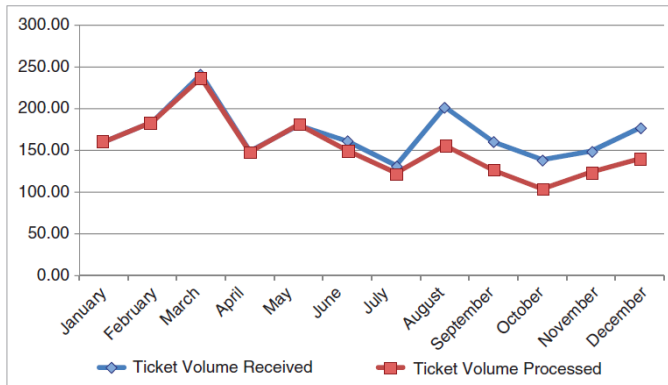


FIGURE 3.17 Original graph

THIS ONE?

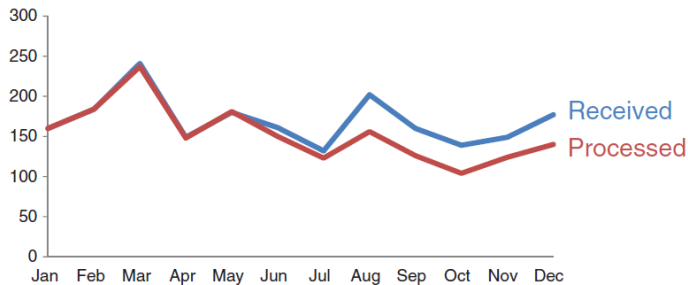


FIGURE 3.23 Leverage consistent color

Top 10 design concerns

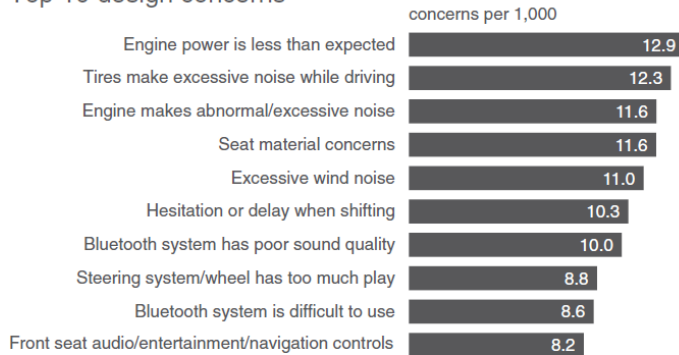


FIGURE 4.7 Original graph, no preattentive attributes

THIS ONE?



Of the top design concerns, three are noise-related.

Top 10 design concerns

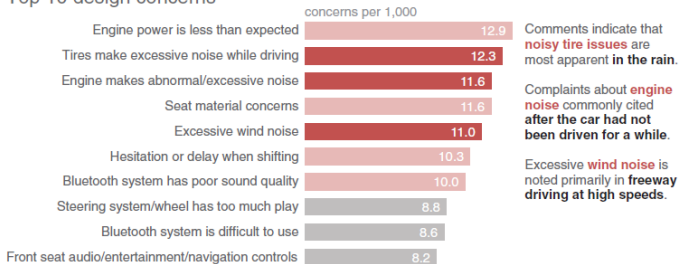
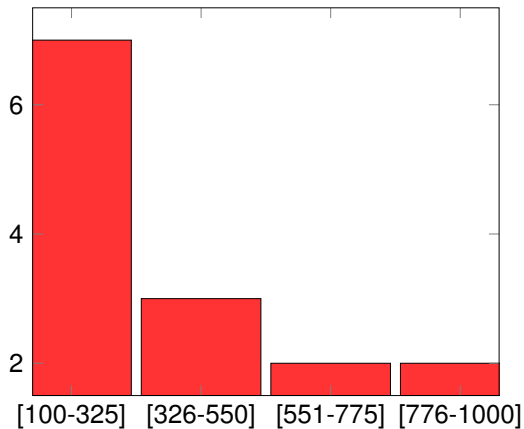


FIGURE 4.9 Create a visual hierarchy of information

EXERCISE

How can we make this better?



- Storytelling with Data by Cole Nussbaumer Knaflic; Wiley

THANK YOU