

ISA 401: Business Intelligence & Data Visualization

14: Fundamentals of Data Visualization

Fadel M. Megahed, PhD

Endres Associate Professor
Farmer School of Business
Miami University

 @FadelMegahed

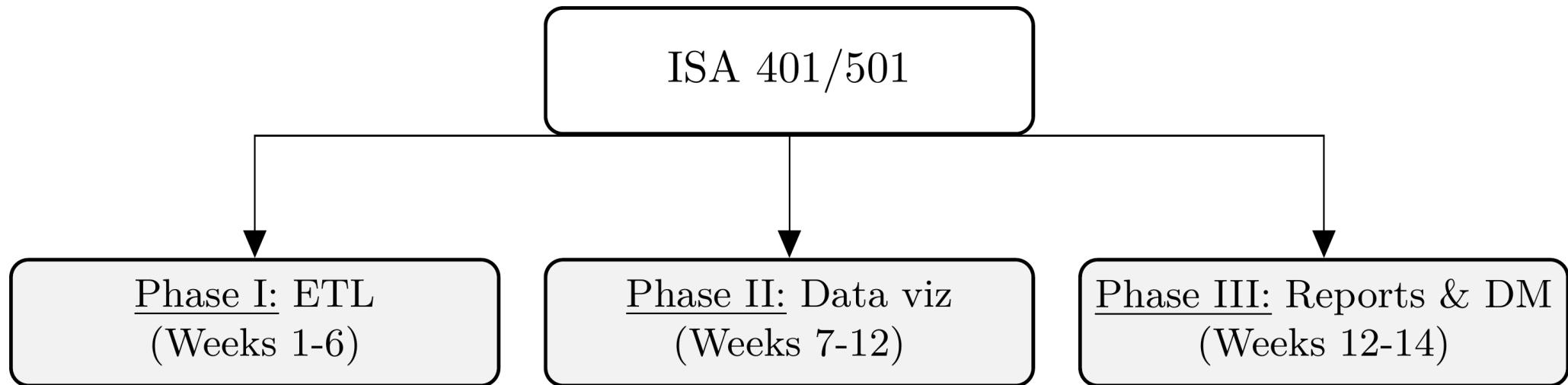
 fmegahed

 fmegahed@miamioh.edu

 Automated Scheduler for Office Hours

Fall 2023

Refresher: Organization of this Course



How the ISA 401/501 course is organized.

Learning Objectives for Today's Class

- Explain the concept of "graphical excellence"
- Explain the theory of data graphics
- Optimize visual encoding based on data types
- Understand why color should be used sparingly and how to select appropriate colors (when color is a must)

Graphical Excellence

Non-graded activity: Terrible Charts

Activity

Russia's Defense Budget

White House Economy Growth

Tucker Carlson

Over the next 5 minutes, please identify the **1-2 main problems** in the charts in the following tabs.

- Write down your answers in the editable area of each chart.
- Discuss your answers with your neighboring classmates.
- Be prepared to share these answers with class.

Non-graded activity: Terrible Charts

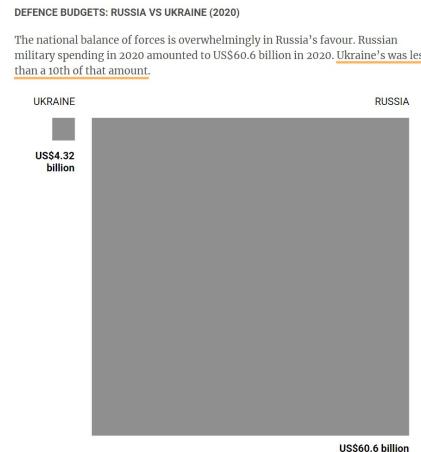
Activity

Russia's Defense Budget

White House Economy Growth

Tucker Carlson

Source: The chart was embedded in [this tweet](#) by Cedric Scherer; however, it is unclear which news outlet have created the original chart.



Main Issue(s): (Insert below)

Non-graded activity: Terrible Charts

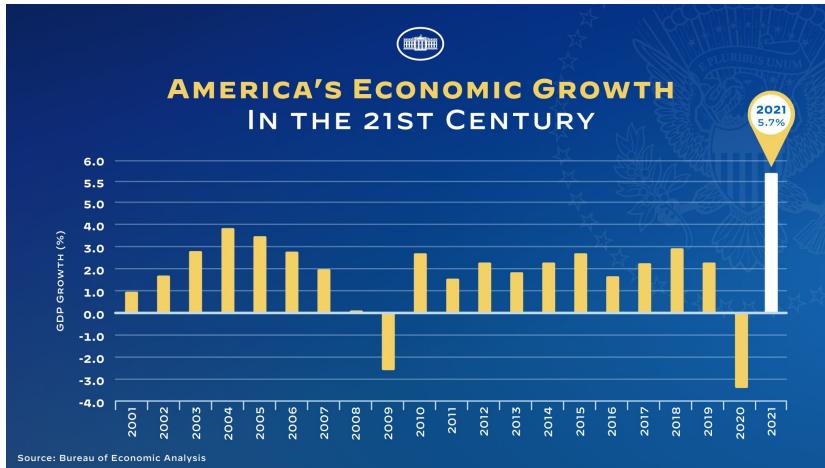
Activity

Russia's Defense Budget

White House Economy Growth

Tucker Carlson

Source: The chart was created by the White House and shared via [this tweet from the verified White House account](#). Note that the chart was latter corrected.



Main Issue(s): (Insert below)

Non-graded activity: Terrible Charts

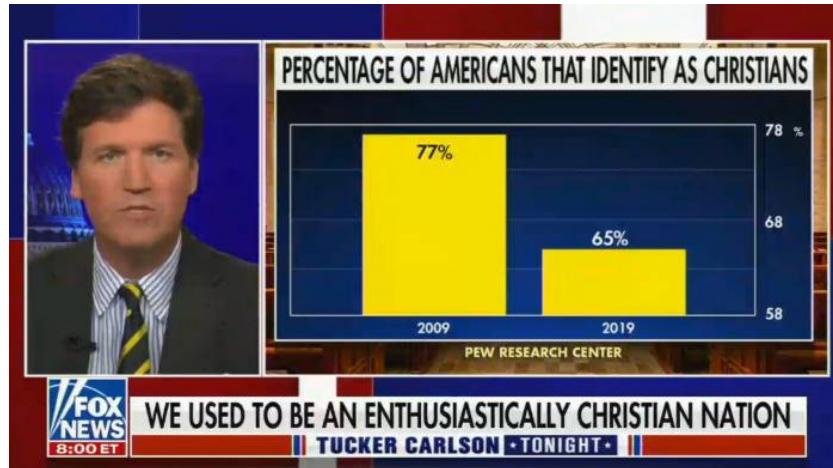
Activity

Russia's Defense Budget

White House Economy Growth

Tucker Carlson

Source: The chart was created by Fox and was highlighted in [this Fox News Clip](#).



Main Issue(s): (Insert below)

Graphical Excellence: What should Graphs Do?

- Show the data
- Lead to thinking about the **substance** rather than something else
- Avoid **distorting** what the data have to say
- Present **many numbers in a small space**
- Make **large datasets coherent**
- Encourage the eye to **compare different pieces of the data**
- **Reveal the data at several levels of detail**, from a broad overview to the fine structure
- Serve **a purpose**: description, exploration, tabulation, decoration
- Be **closely integrated with the statistical & verbal descriptions of the data**

Show/Reveal the Data: Anscombe's Dataset

In a seminal paper, Anscombe stated:

Few of us escape being indoctrinated with these notions:

- numerical **calculations are exact, but graphs are rough**;
- for any particular kind of **statistical data there is just one set of calculations constituting a correct statistical analysis**;
- performing **intricate calculations is virtuous**, whereas **actually looking at the data is cheating**.

He proceeded by stating that

a computer should **make both calculations and graphs**. Both sorts of output should be studied; each will contribute to understanding.

Now, let us consider his four datasets, each consisting of eleven (x,y) pairs.

.footnote[

Source: Anscombe, Francis J. 1973. "Graphs in Statistical Analysis." *The American Statistician* 27 (1): 17–21. ([PDF Link](#)).

Show/Reveal the Data: Anscombe's Dataset

x1 ◆	x2 ◆	x3 ◆	x4 ◆	y1 ◆	y2 ◆	y3 ◆	y4 ◆
10	10	10	8	8.04	9.14	7.46	6.58
8	8	8	8	6.95	8.14	6.77	5.76
13	13	13	8	7.58	8.74	12.74	7.71
9	9	9	8	8.81	8.77	7.11	8.84
11	11	11	8	8.33	9.26	7.81	8.47
14	14	14	8	9.96	8.1	8.84	7.04
6	6	6	8	7.24	6.13	6.08	5.25
4	4	4	19	4.26	3.1	5.39	12.5
12	12	12	8	10.84	9.13	8.15	5.56
7	7	7	8	4.82	7.26	6.42	7.91
5	5	5	8	5.68	4.74	5.73	6.89

Showing 1 to 11 of 11 entries

Previous 1 Next

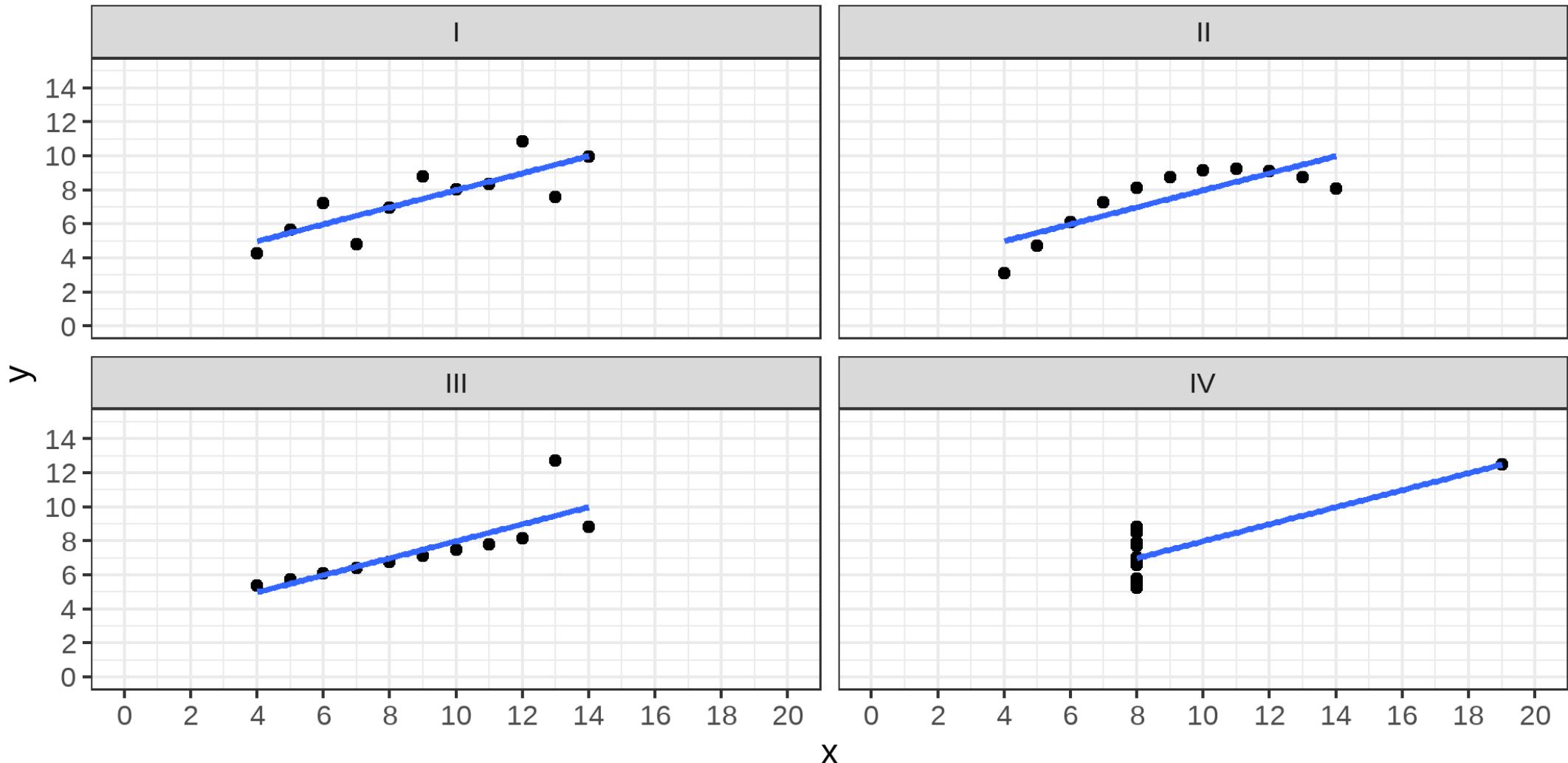
Show/Reveal the Data: Anscombe's Dataset

set	x.mean	x.sd	y.mean	y.sd	corr
I	9	3.32	7.5	2.03	0.82
II	9	3.32	7.5	2.03	0.82
III	9	3.32	7.5	2.03	0.82
IV	9	3.32	7.5	2.03	0.82

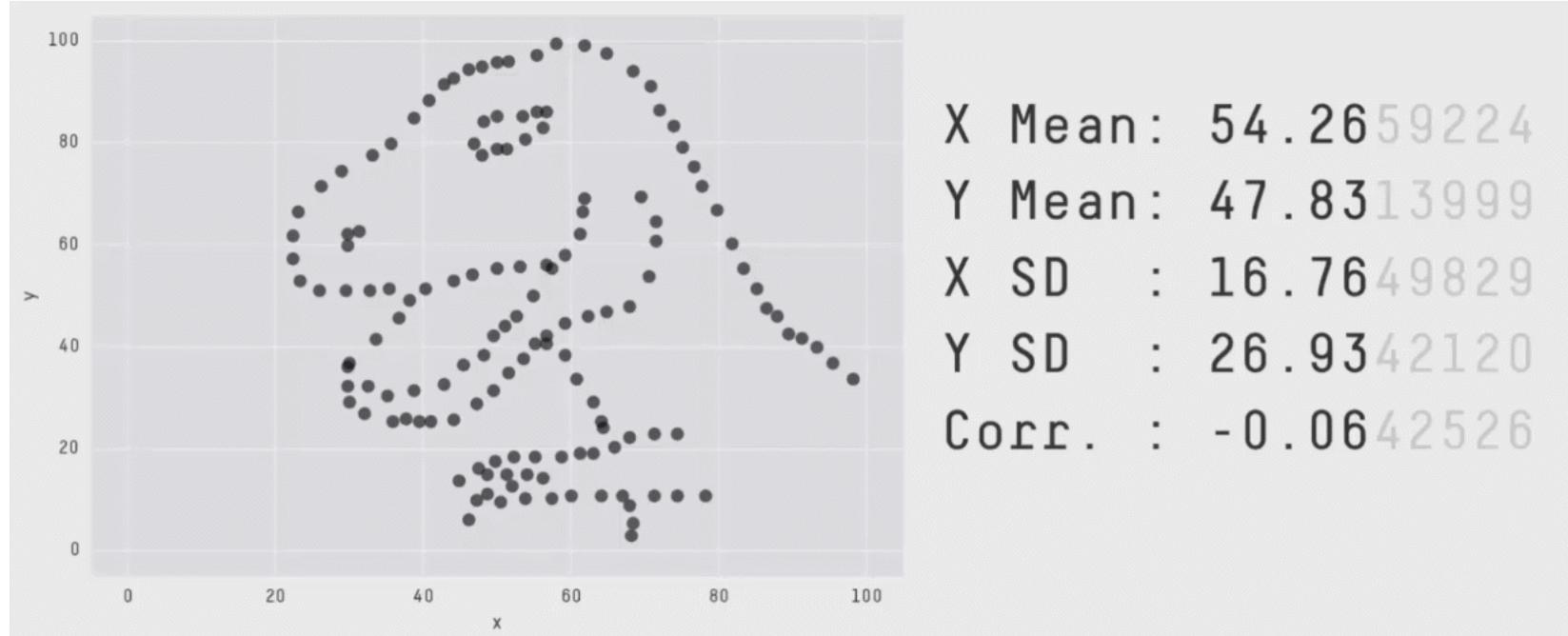
Showing 1 to 4 of 4 entries

Previous 1 Next

Show/Reveal the Data: Anscombe's Dataset



A Modern Version of Anscombe's Dataset



.footnote[

Source: Matejka, J. and Fitzmauricem G. 2023. "Same Stats, Different Graphs" *Proceedings of the 2017 CHI conference on human factors in computing systems.* ([Blog Post Link](#)).

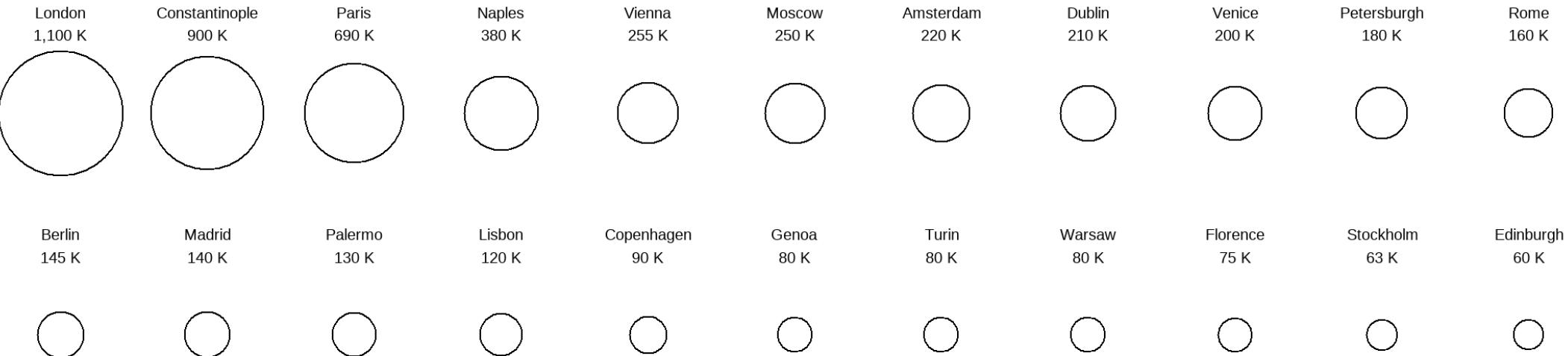
Substance

Activity

Your Solution

My Solution

In 5 minutes, please **sketch** a better (non-bubble) chart than the one use by William Playfair for plotting the populations of 22 European cities at the end of the 1700s.



Source: Chart made by Fadel M. Megahed based on data from <http://www.stat.uiowa.edu/~luke/data/Playfair>

Substance

Activity

Your Solution

My Solution

Ideally, on a piece of paper sketch out your solution. Otherwise, please feel free to download the plot's data (using the code below) and use a software of your choice for plotting a better chart for the data.

```
pacman::p_load(tidyverse)
playfair = read.table("http://www.stat.uiowa.edu/~luke/data/Playfair") %>%
  rownames_to_column(var = 'city') %>% # converting row names to city var
  as_tibble() %>% # converting it to a tibble
  arrange( desc(population) ) # arranging the rows in a descending order by population
write_csv(x = playfair, file = 'playfair_data.csv')
```

Be prepared to share your solution with the entire class.

Substance

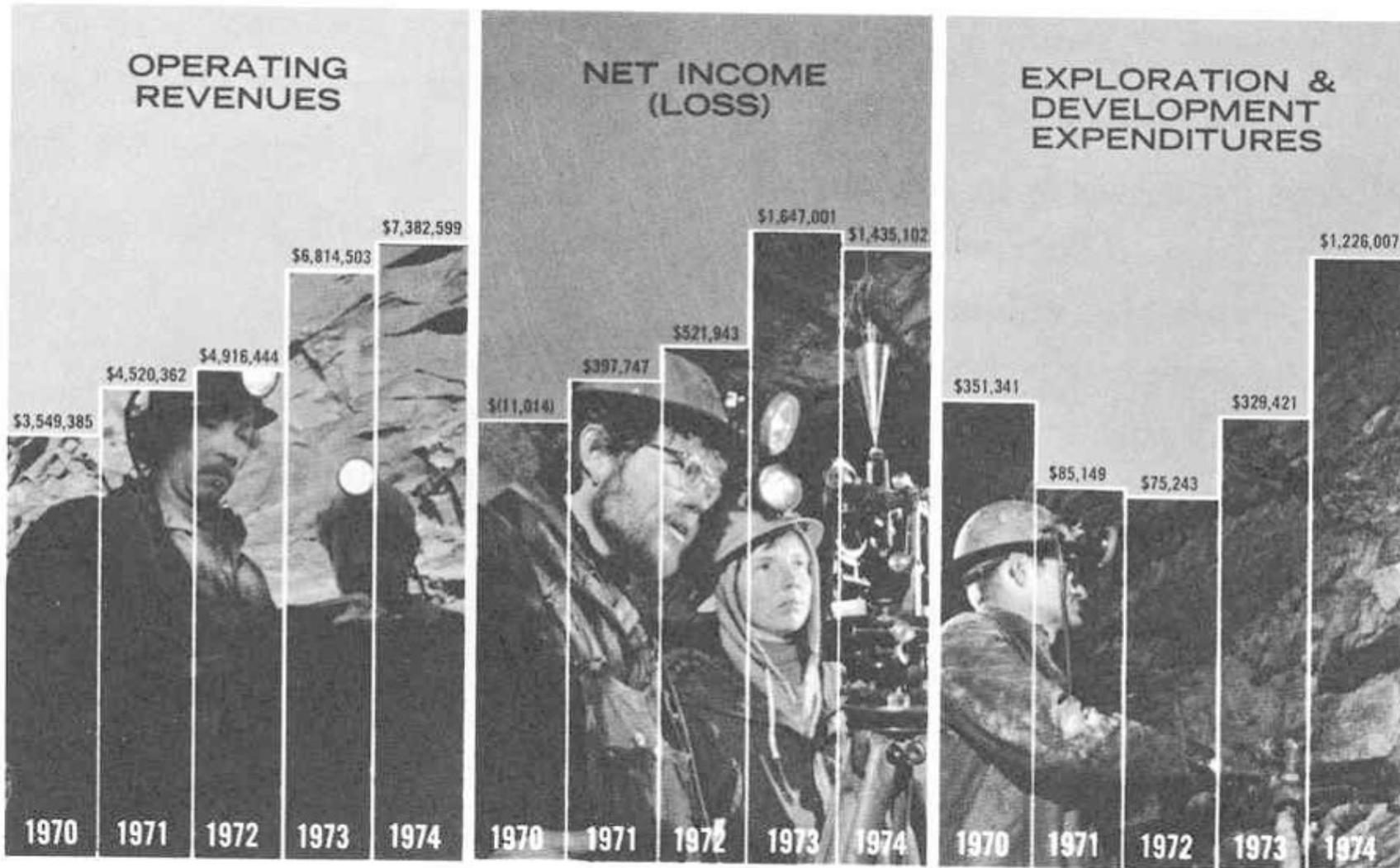
Activity

Your Solution

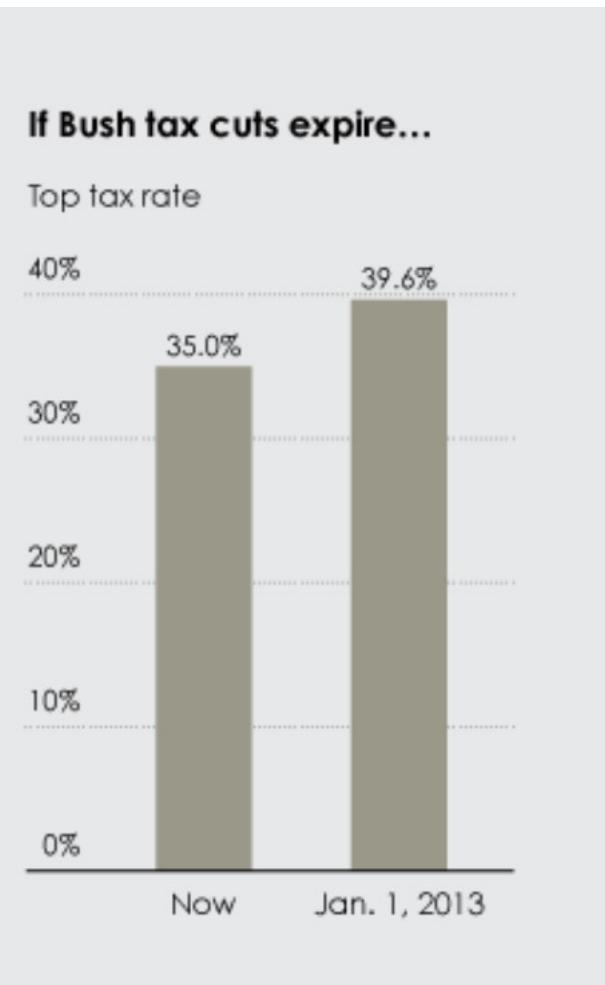
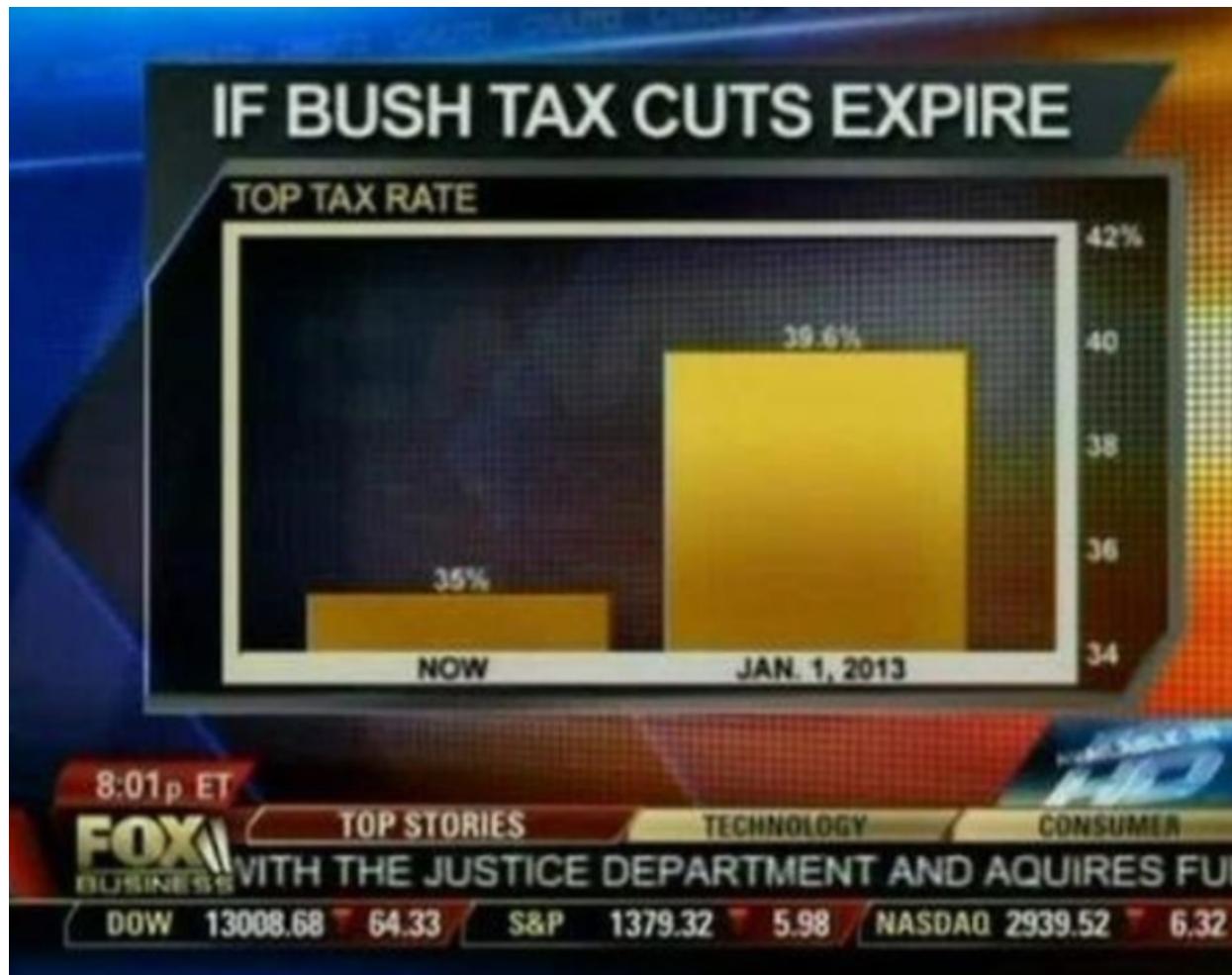
My Solution

- In my opinion, a **dot chart is more effective** than the bubble chart. The **population would be mainly encoded using the position**; you can still use area as a secondary encoding mechanism.

Avoid Distortion of Data



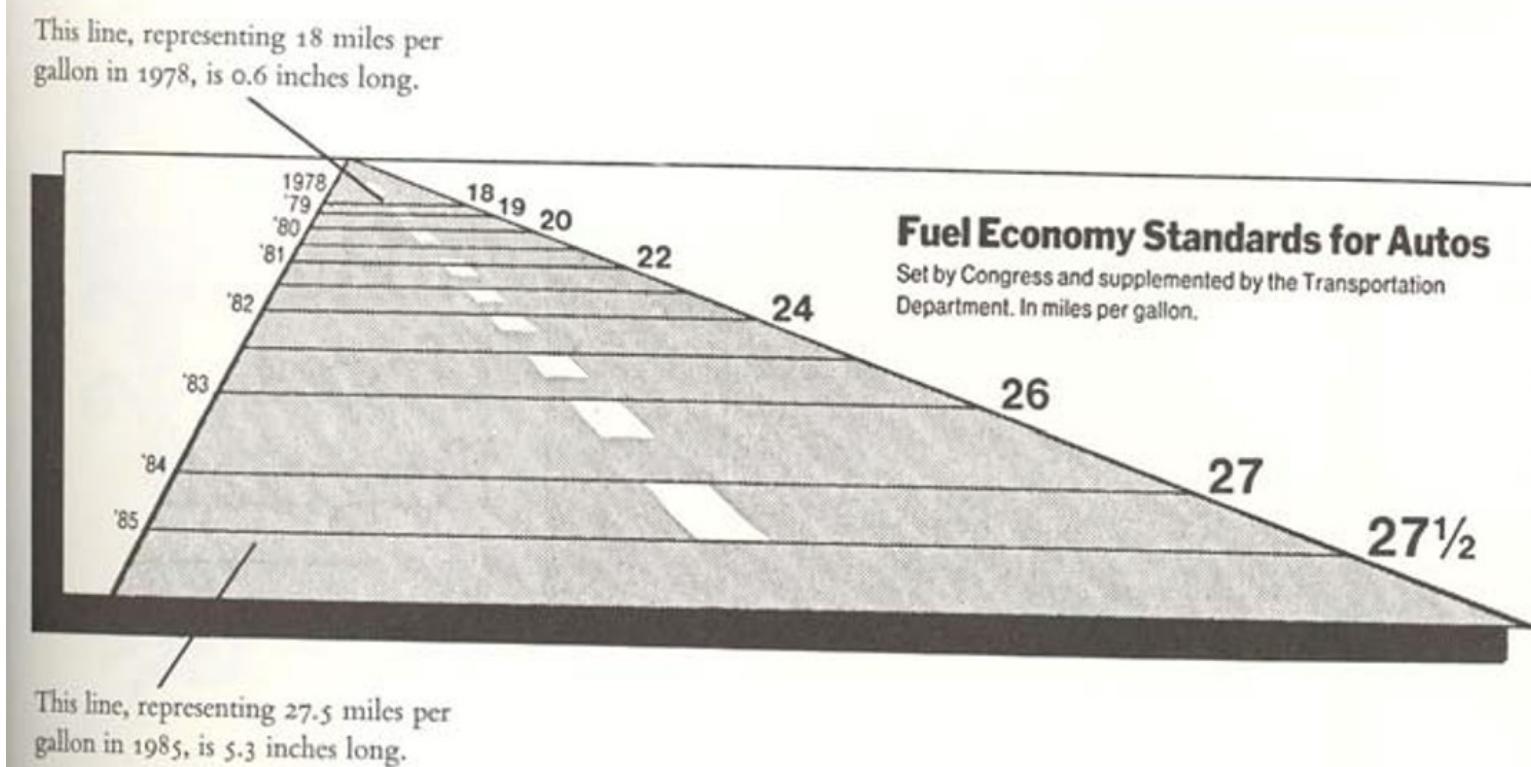
Avoid Distortion of Data



Avoid Distortion of Data: The Lie Factor

Size of effect shown in graphic

Size of effect in data



Graphical Integrity Principles: A Summary

- Clear, detailed, and thorough labeling and appropriate scales
- Size of the graphic effect should be directly proportional to the numerical quantities ("lie factor")
- Show data variation, not design variation

Theory of Data Graphics

Definition of Data Ink

- **Data-ink refers to the non-erasable ink used for presenting the data.**
 - If data-ink would be removed from the image, the graphic would lose the content.

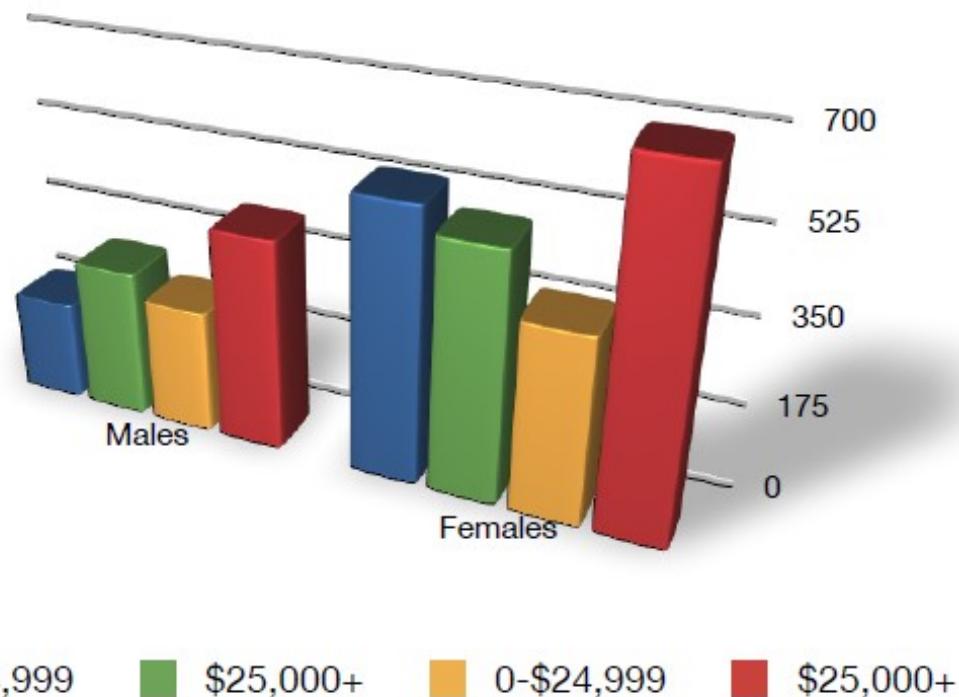
$$\text{Data-ink ratio} = \frac{\text{Data-ink}}{\text{Total ink used to print the graphic}}$$

= proportion of a graphic's ink devoted to the
non-redundant display of data-information

= $1.0 - \text{proportion of a graphic that can be erased}$

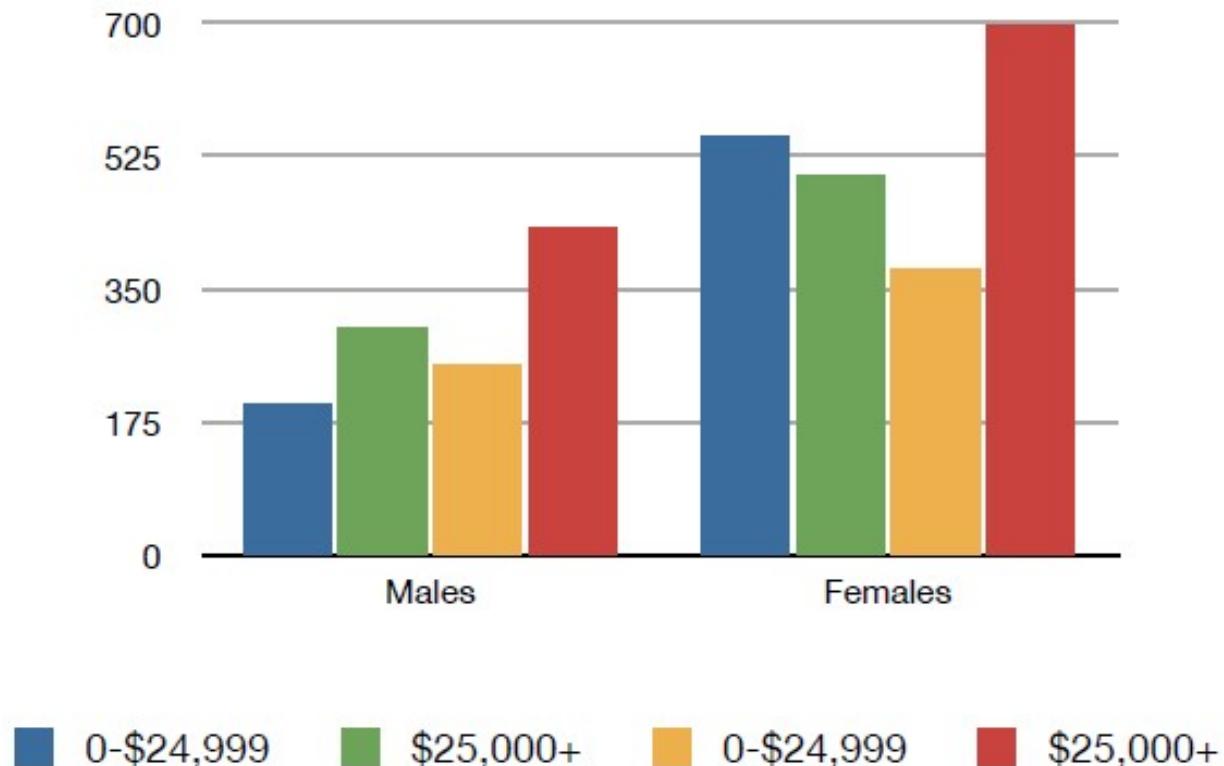
Definition of Data Ink

- Data-ink refers to the non-erasable ink used for presenting the data.
 - If data-ink would be removed from the image, the graphic would lose the content.



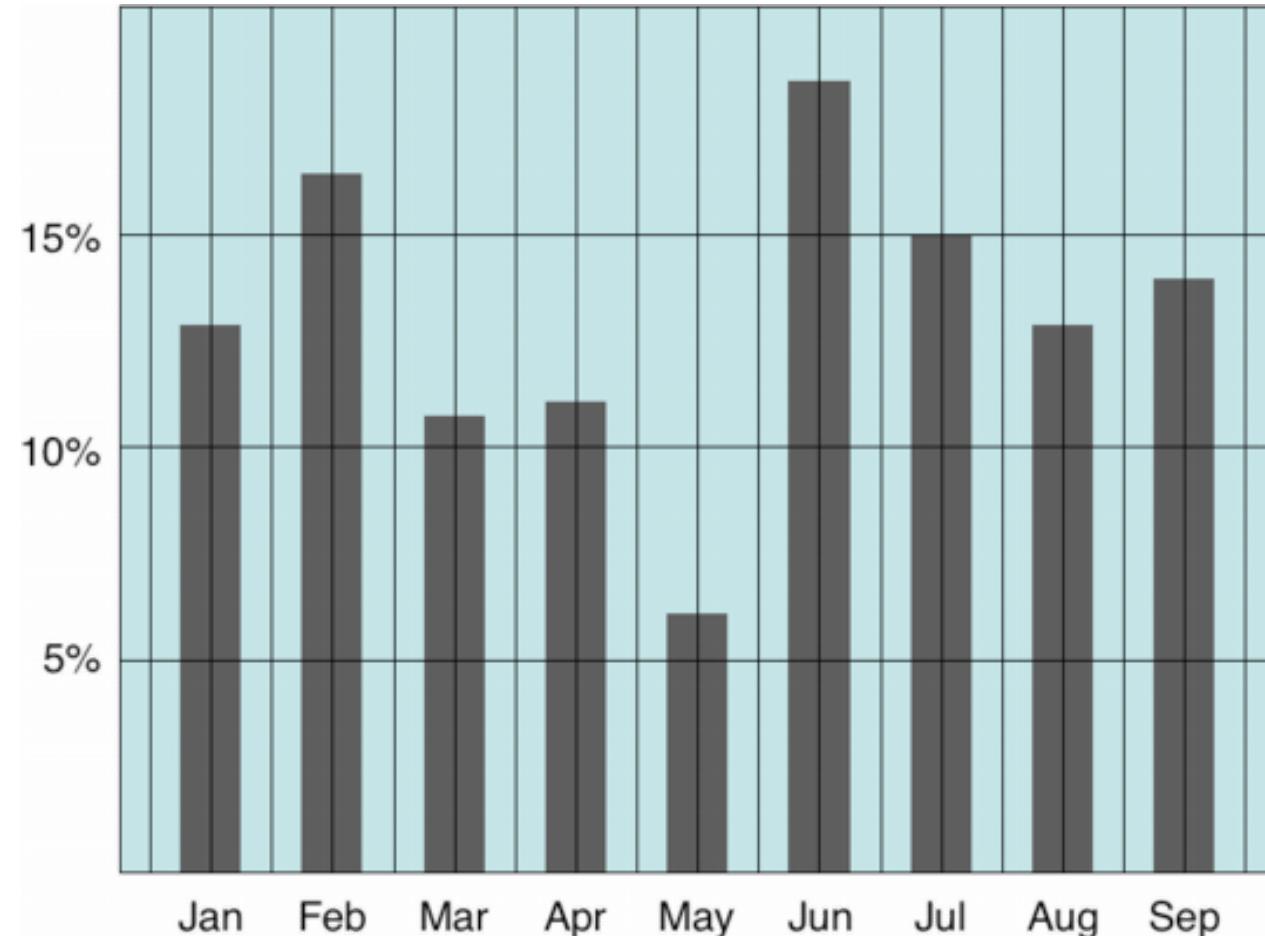
Definition of Data Ink

- Data-ink refers to the non-erasable ink used for presenting the data.
 - If data-ink would be removed from the image, the graphic would lose the content.



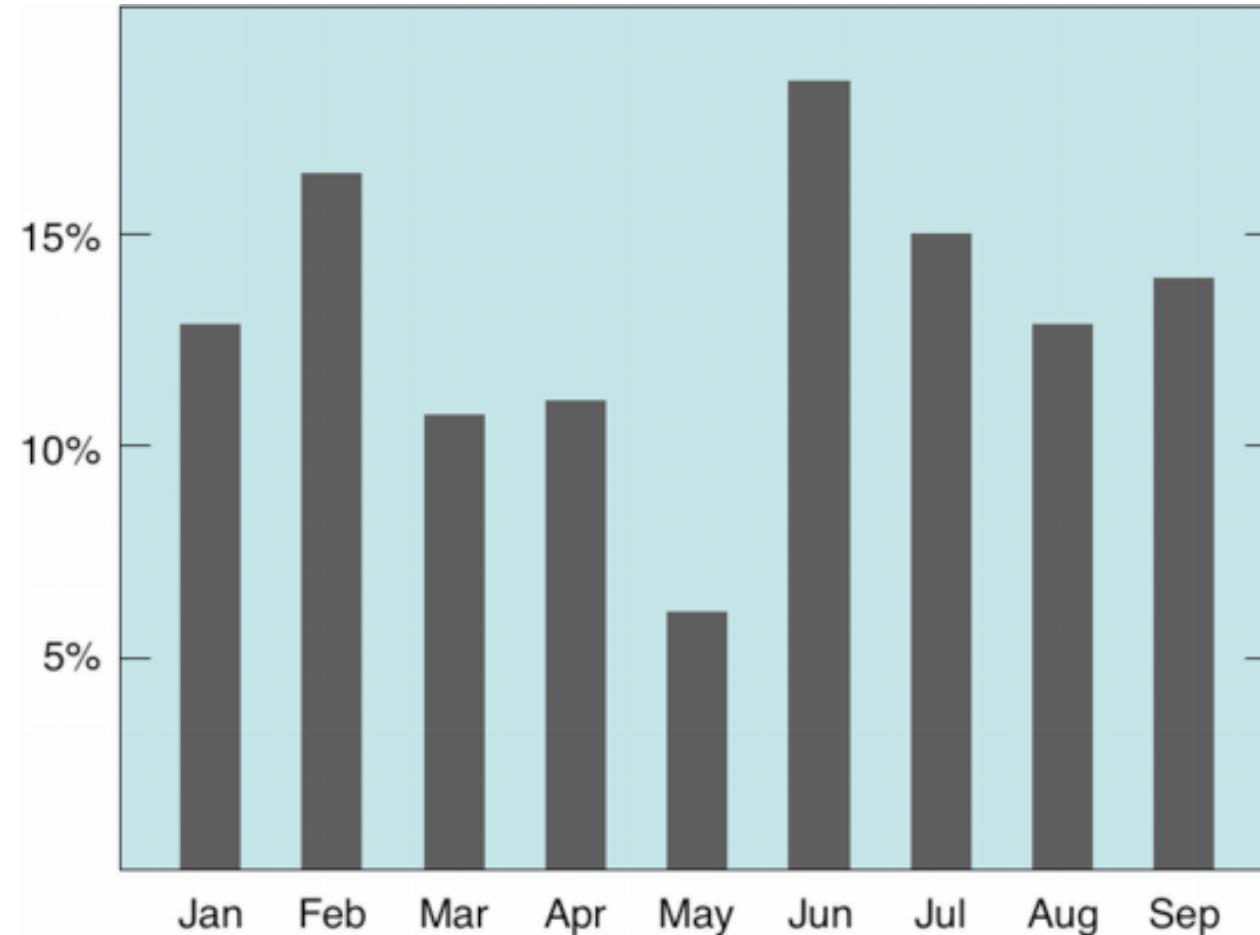
Focus on Data: Avoid Chartjunk

Chartjunk is the extraneous visual elements that distract from the message!!



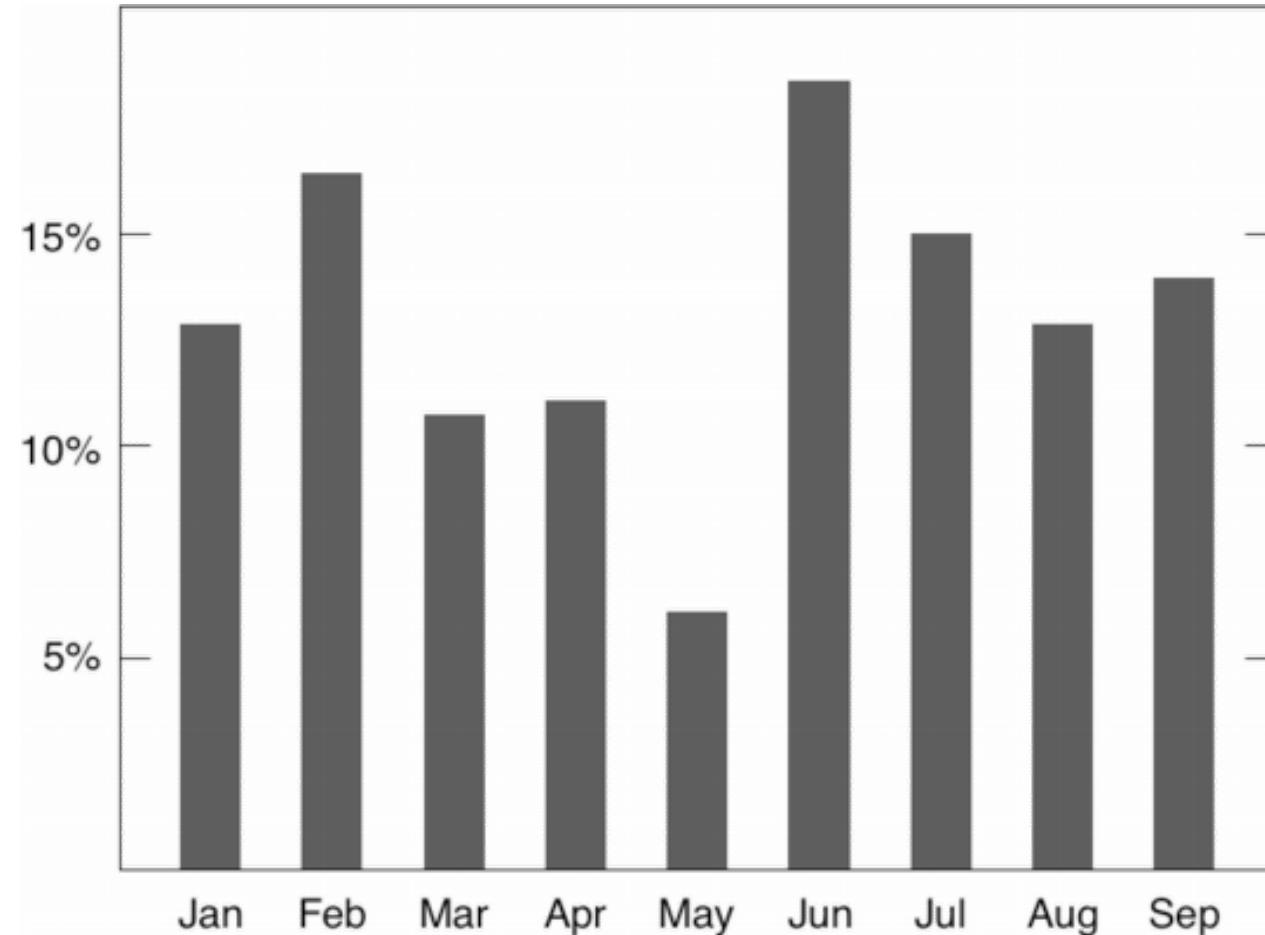
Focus on Data: Avoid Chartjunk

Chartjunk is the extraneous visual elements that distract from the message!!



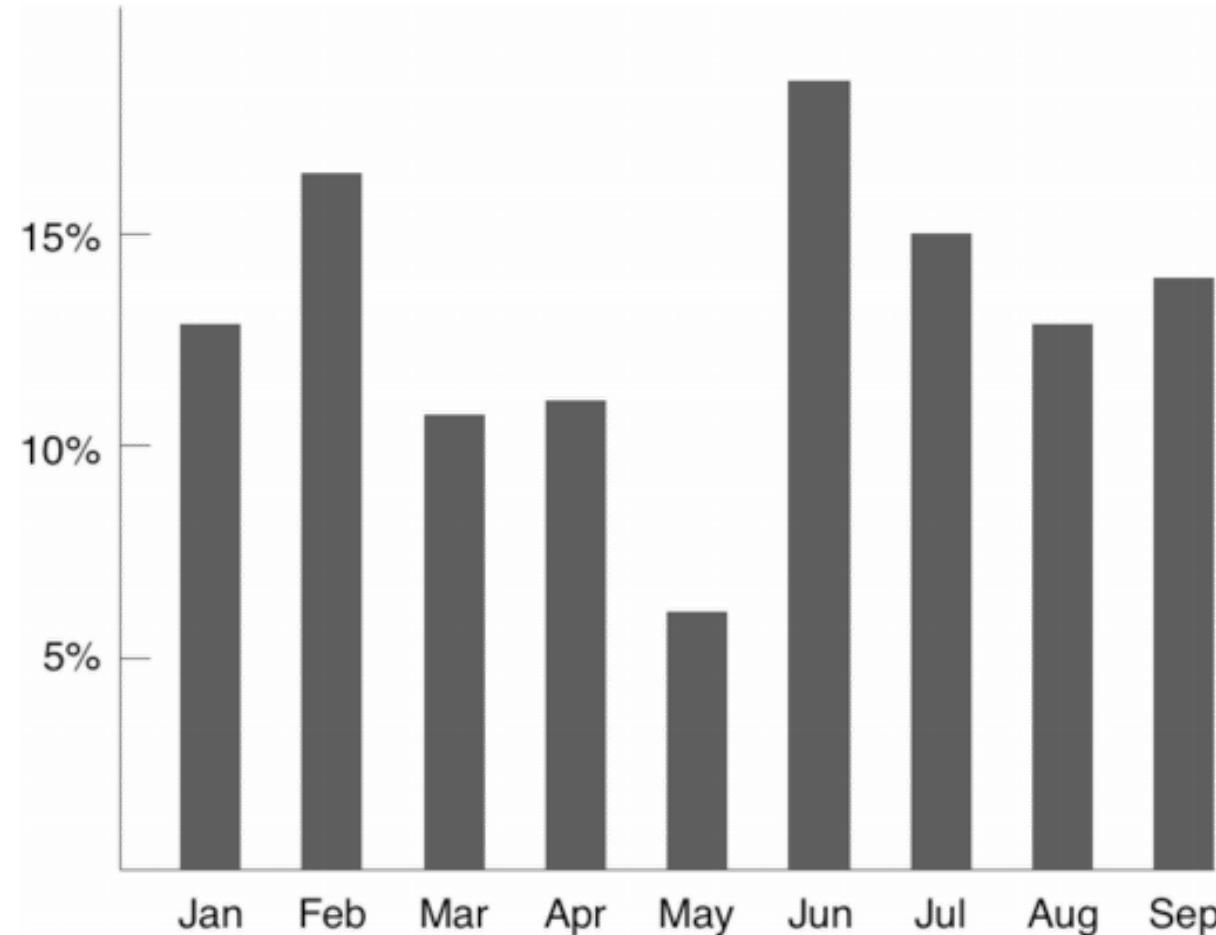
Focus on Data: Avoid Chartjunk

Chartjunk is the extraneous visual elements that distract from the message!!



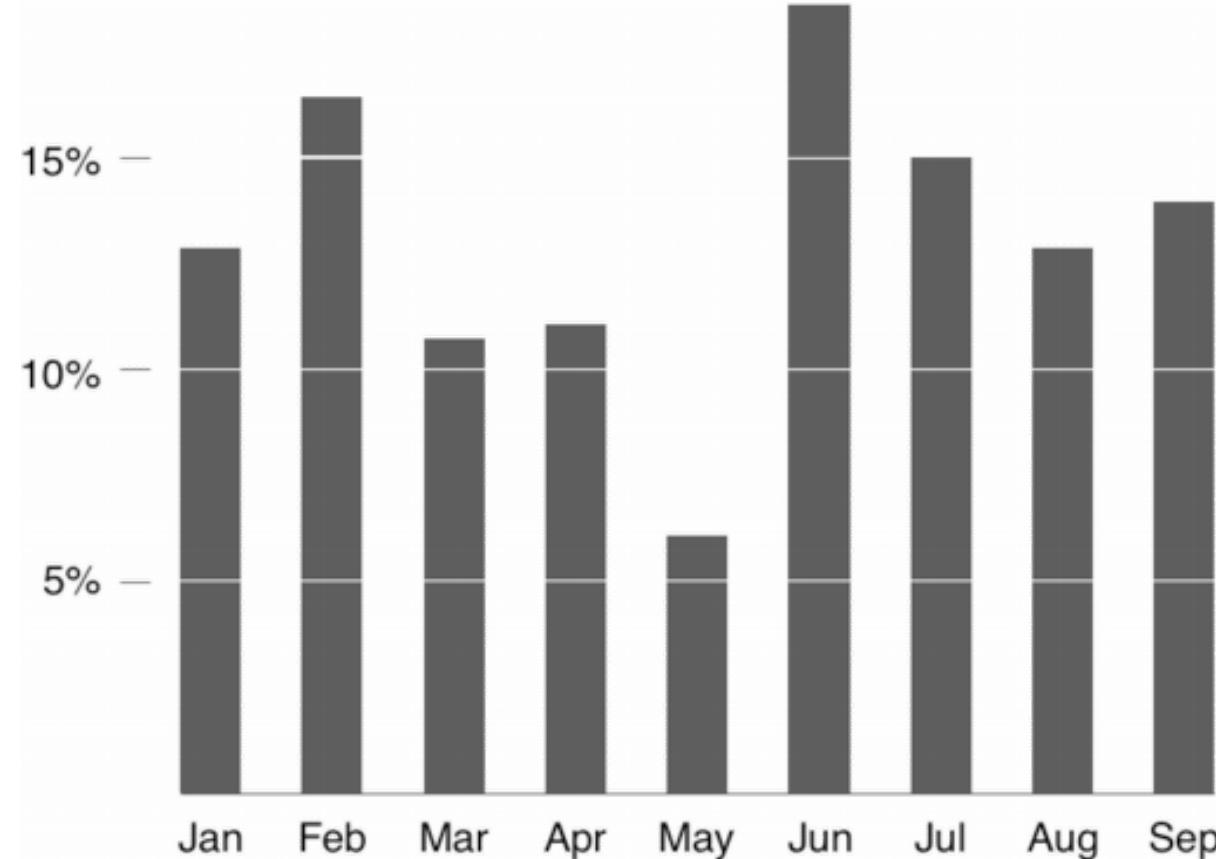
Focus on Data: Avoid Chartjunk

Chartjunk is the extraneous visual elements that distract from the message!!



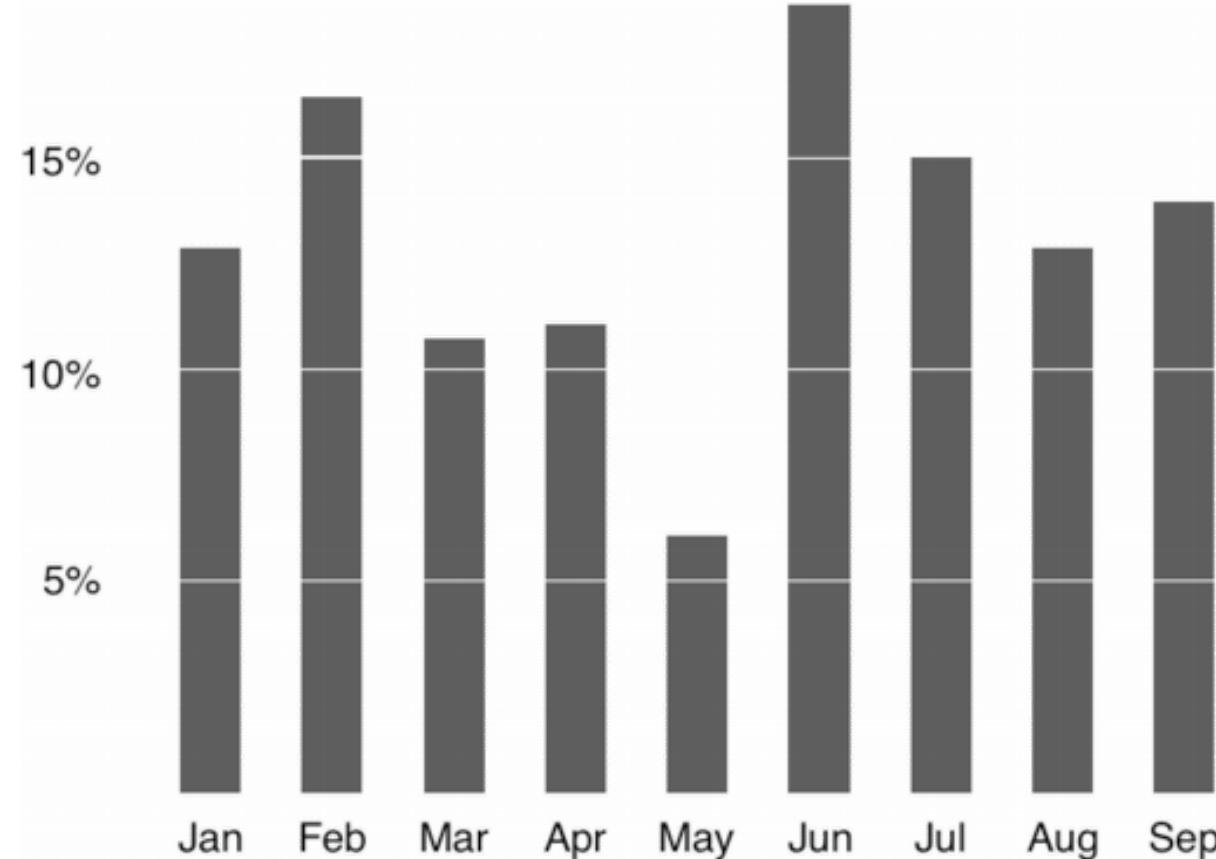
Focus on Data: Avoid Chartjunk

Chartjunk is the extraneous visual elements that distract from the message!!



Focus on Data: Avoid Chartjunk

Chartjunk is the extraneous visual elements that distract from the message!!



Other Subjective Design Principles

- **Aesthetics:** Attractive things are perceived as more useful than unattractive ones
- **Style:** Communicates brand, process, who the designer is
- **Playfulness:** Encourages experimentation and exploration
- **Vividness:** Can make a visualization more memorable

Data Models

SCIENCE

Vol. 103, No. 2684

Friday, June 7, 1946

On the Theory of Scales of Measurement

S. S. Stevens

Director, Psycho-Acoustic Laboratory, Harvard University

FOR SEVEN YEARS A COMMITTEE of the British Association for the Advancement of Science debated the problem of measurement. Appointed in 1932 to represent Section A (Mathematical and Physical Sciences) and Section J (Psychology), the committee was instructed to consider and report upon the possibility of "quantitative estimates of sensory events"—meaning simply: Is it possible to measure human sensation? Deliberation led only to disagreement, mainly about what is meant by the term measurement. An interim report in 1938 found one member complaining that his colleagues "came out by that same door as they went in," and in order to have another try at agreement, the committee begged to be continued for another year.

For its final report (1940) the committee chose a common bone for its contentions, directing its arguments at a concrete example of a sensory scale. This was the Sone scale of loudness (S. S. Stevens and H. Davis. *Hearing*. New York: Wiley, 1938), which purports to measure the subjective magnitude of an auditory sensation against a scale having the formal properties of other basic scales, such as those used to measure length and weight. Again the 19 members of the committee came out by the routes they entered, and their views ranged widely between two extremes. One member submitted "that any law purporting to express a quantitative relation between sensation intensity and stimulus intensity is not merely false but is in fact meaningless unless and until a meaning can be given to the concept of addition as applied to sensation" (Final Report, p. 245).

It is plain from this and from other statements by the committee that the real issue is the meaning of measurement. This, to be sure, is a semantic issue, but one susceptible of analysis. Perhaps

by the formal (mathematical) properties of the scales. Furthermore—and this is of great concern to several of the sciences—the statistical manipulations that can legitimately be applied to empirical data depend upon the type of scale against which the data are ordered.

A CLASSIFICATION OF SCALES OF MEASUREMENT

Paraphrasing N. R. Campbell (Final Report, p. 340), we may say that measurement, in the broadest sense, is defined as the assignment of numerals to objects or events according to rules. The fact that numerals can be assigned under different rules leads to different kinds of scales and different kinds of measurement. The problem then becomes that of making explicit (a) the various rules for the assignment of numerals, (b) the mathematical properties (or group structure) of the resulting scales, and (c) the statistical operations applicable to measurements made with each type of scale.

Scales are possible in the first place only because there is a certain isomorphism between what we can do with the aspects of objects and the properties of the numeral series. In dealing with the aspects of objects we invoke empirical operations for determining equality (classifying), for rank-ordering, and for determining when differences and when ratios between the aspects of objects are equal. The conventional series of numerals yields to analogous operations: We can identify the members of a numeral series and classify them. We know their order as given by convention. We can determine equal differences, as $8 - 6 = 4 - 2$, and equal ratios, as $8/4 = 6/3$. The isomorphism between these properties of the numeral series and certain empirical operations which we perform with objects permits the use of the series as a model to represent aspects of the empirical world.

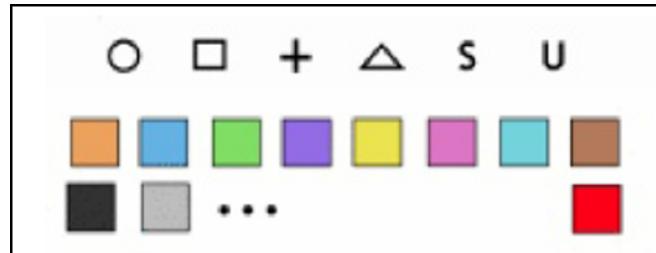
Data Types (from S. Stevens, Theory of Scales)

Scale	Basic Empirical Operations	Mathematical Group Structure	Permissible Statistics (invariantive)
NOMINAL	Determination of equality	<i>Permutation group</i> $x' = f(x)$ $f(x)$ means any one-to-one substitution	Number of cases Mode Contingency correlation
ORDINAL	Determination of greater or less	<i>Isotonic group</i> $x' = f(x)$ $f(x)$ means any monotonic increasing function	Median Percentiles
INTERVAL	Determination of equality of intervals or differences	<i>General linear group</i> $x' = ax + b$	Mean Standard deviation Rank-order correlation Product-moment correlation
RATIO	Determination of equality of ratios	<i>Similarity group</i> $x' = ax$	Coefficient of variation

Data Types: Explained

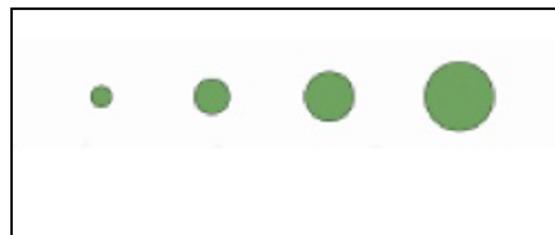
Nominal:

- Are = or \neq to other values
- Apples, bananas, oranges, etc.



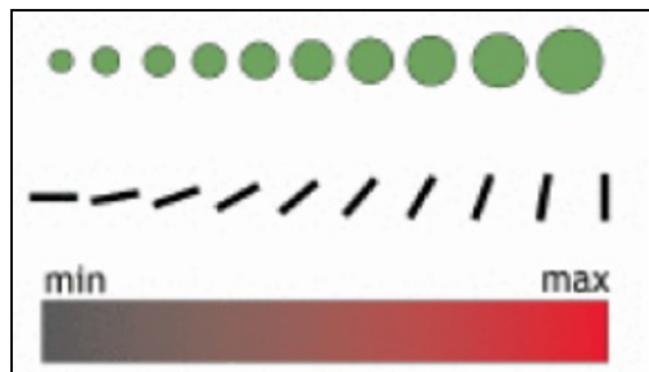
Ordinal:

- Obey a $<$ relationship
- Small, medium and large



Quantitative:

- Can do math on them
- 50 inches, 53 inches, etc.



Quantitative Data Types (from S. Stevens)

Quantitative data can be further divided into:

Intervals (Location of Zero Arbitrary):

- Dates: Jan 19; Location: (Lat, Long)
- Only differences (i.e., intervals) can be compared.

Ratio (Zero Fixed):

- Measurements: Length, Weight, ...
- Origin is meaningful, we can compute ratios, proportions, differences, etc.

Non-Graded Activity: Data Terminology

Activity Your Solution

In 2 minutes, please identify an appropriate data type for each column below.

Order ID	Order Date	Order Priority	Product Container	Product Cost	Ship Date
1	1/1/2022	5 - low	Large box	25	1/5/2022
2	1/4/2022	4 - not specified	Small Box	36	1/7/2022
3	1/15/2022	2- high	Small Box	38	1/17/2022
3	1/15/2022	2- high	Small Box	41	1/17/2022
3	1/15/2022	2- high	Jumbo Box	44	1/17/2022
3	1/15/2022	2- high	Wrap Bag	33	1/17/2022
4	1/18/2022	1- urgent	Small Box	33	1/19/2022

Showing 1 to 7 of 11 entries

Previous 1 2 Next

02:00

Non-Graded Activity: Data Terminology

Activity Your Solution

Data Types: (Edit below)

- Order ID:
- Order/Ship Date:
- Order Priority:
- Product Container:
- Product Cost:

Data vs. Conceptual Models

From data model

- 32.5, 54.0, -17.3, ...

Using a conceptual model:

- Temperature

To data type:

- Continuous to x significant digits i.e. quantitative
- Hot, warm, cold i.e. ordinal
- Burned vs. not burned i.e. nominal

Image Model: Visual (Encoding) Variables

Channels	Marks	Points	Lines	Areas
Position	XY 2 DIMENSIONS DU PLAN	POINTS	LIGNES	ZONES
Size	Z TAILLE			
(Grey) Value	VALEUR			
Texture	GRAIN			
Color	COULEUR			
Orientation	ORIENTATION			
Shape	FORME			

LES VARIABLES DE L'IMAGE

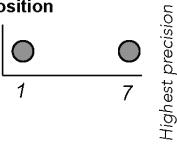
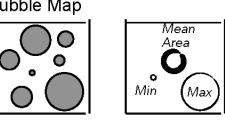
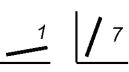
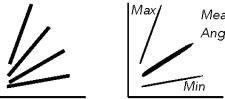
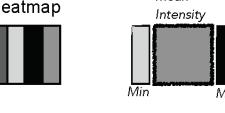
LES VARIABLES DE SÉPARATION DES IMAGES

Mapping to Data Types

	Nominal	Ordinal	Quantitative
Position	✓	✓	✓
Size	✓	✓	~
(Grey)Value	✓	✓	~
Texture	✓	~	✗
Color	✓	✗	✗
Orientation	✓	✗	✗
Shape	✓	✗	✗

✓ = Good
~ = OK
✗ = Bad

Visual Channels and their Precision

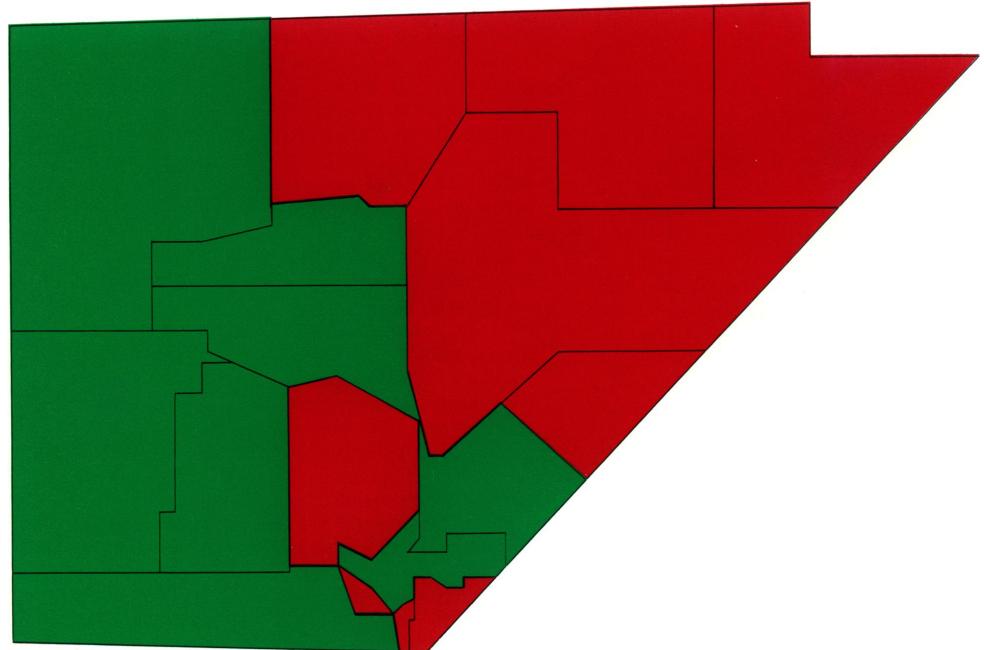
Absolute precision ranking for seeing a single ratio	Common illusions that distort data	Vision is powerful for global statistics	Vision is sluggish for comparisons
Visual encodings and their differences in precision for estimating ratios (here, 1:7)	Caveats for the visual encoding in each row	For each visualization, statistics are available quickly	Isolating pairs with 'larger second values' is tough
Position 	 Stacked bar: Bars on baseline are position-coded = more precise perception. The black & dark gray bars have the same value differences among them, but the differences are only visible across the black bars.	Dot Plot Stacked Bar Bubble Map Slope Graph Heatmap 	 Tool: Shortcut comparisons by adding direct depictions of the deltas, as below Tool: Highlight and annotate the right comparisons for your viewers, as above Tool: For color heatmaps, depict deltas as blue (+) & red (-)
Length 			
Area 	If numbers map to Area, this shows a ~1:7 ratio. But beware: if numbers map to Length, then grey circles depict ~1:2.5 ratio.		" $a, c, & e$ have increased" Tool: Highlight and annotate the right comparisons for your viewers, as above
Angle 	The difference is larger for the lighter segments compared to the darker ones, right? That's an illusion - the differences are identical.		Tool: You and your viewers will (generally) compare values that: (1) are close together or connected and (2) have similar colors, in that priority order
Intensity 	Intensity values can look different depending their backgrounds. Don't plot intensities on intensities.		For color heatmaps, depict deltas as blue (+) & red (-)

Color Should Be Used Sparingly

Color

Activity Your Solution Results from 1983 Experiment

- The following map of Nevada has been colored to indicate various geological features in each county.
- Estimate the larger land area-more red, more green, or the same-and mark your answer on the mentimeter poll in the next panel.
- **Please work fairly quickly, as if you were trying to gain an overall impression from a map.**



Color

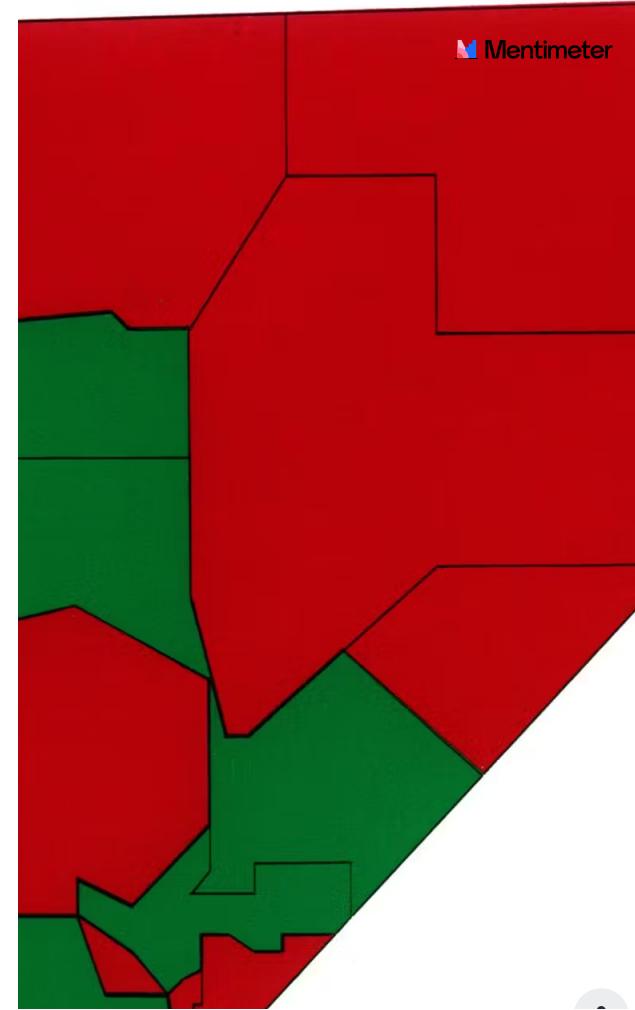
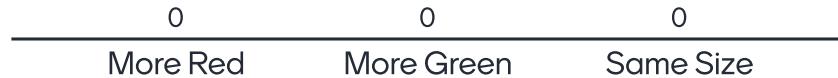
Activity

Your Solution

Results from 1983 Experiment

Join at <https://www.menti.com/g9wh7os8je>

Which is the larger land-area?



Color

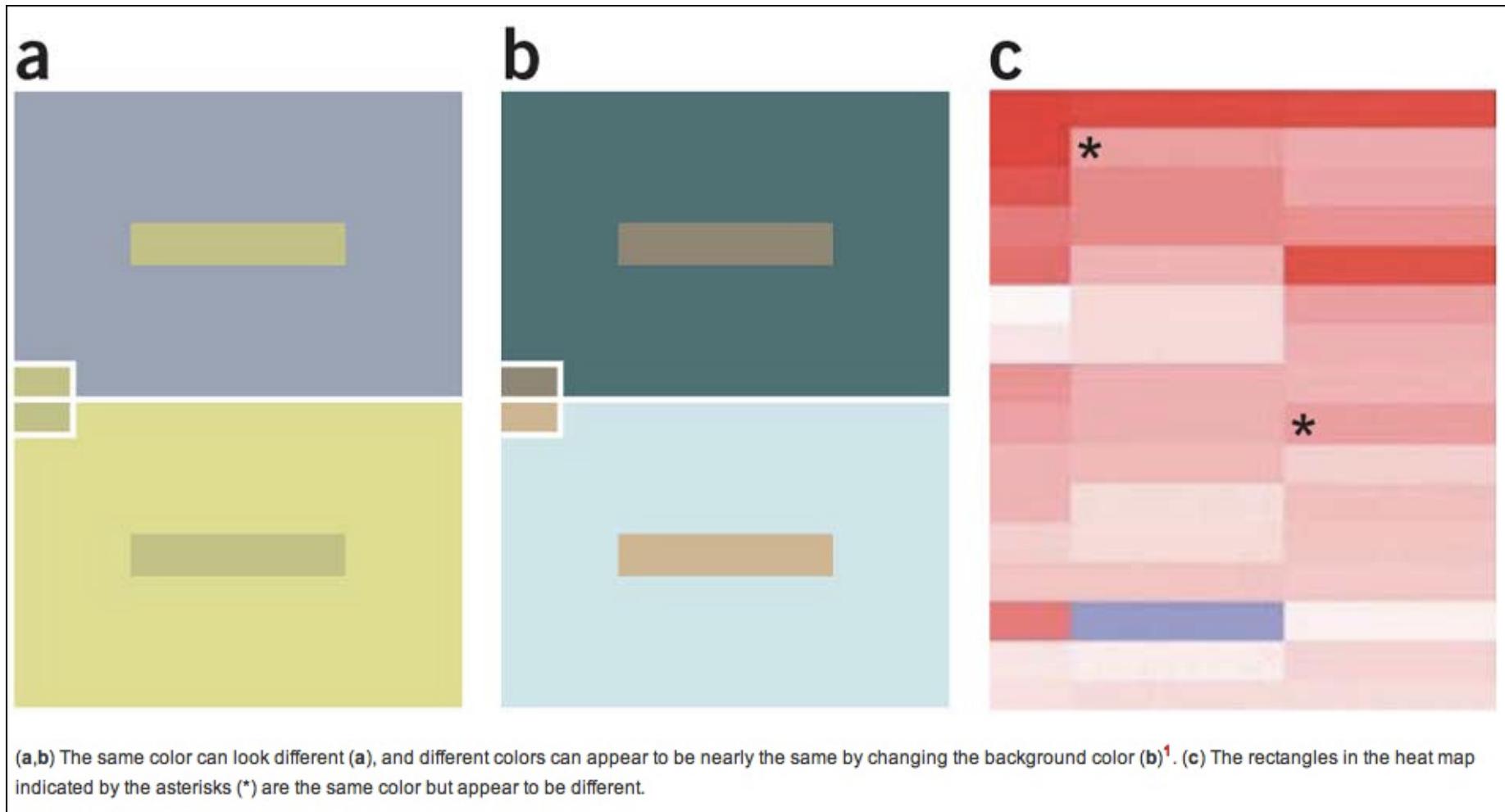
Activity

Your Solution

Results from 1983 Experiment

- The results from the original experiment were as follows

Simultaneous Contrast Affects Perception



Color Blindness

"About 1 in 12 men are color blind" -- NIH's At a glance: Color Blindness

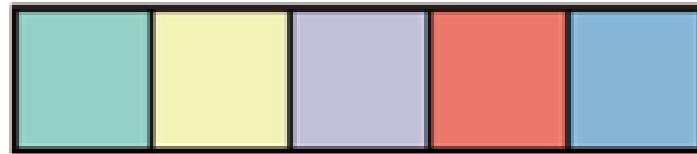
Color Blindness Can Distort a Person's Reading/Interpretation of a Chart

- **Personal Check:** To have a feel for color blindness (if you are not color blind), you can take this color blind test
- Given the high prevalence of color blind individuals, your charts **should** accommodate for color-blindness. **How?**
 - Use color sparingly
 - Use color friendly palettes, e.g., see <https://colorbrewer2.org/>

Color Brewer: Color Scales and their Selection

Nominal

Qualitative Scale

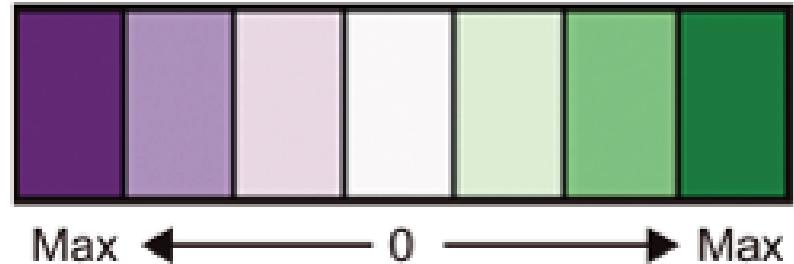


Ordinal

Sequential Scale



Diverging Scale



Recap

Summary of Main Points

By now, you should be able to do the following:

- Explain the concept of "graphical excellence"
- Explain the theory of data graphics
- Optimize visual encoding based on data types
- Understand why color should be used sparingly and how to select appropriate colors (when color is a must)

Things to Do Prior to Next Class

Please go through the following two supplementary readings and complete [assignment 09](#).

- [The Lie Factor and the Baseline Paradox](#); especially noting what the authors mean by "baseline", how the lie factor may be ignored in time-series applications, and/or in applications involving a "ratio" scale.
- [Useful junk? The effects of visual embellishment on comprehension and memorability of charts](#), which presents an experimental counter against Tufte's argument for simplicity (by quantifying vividness and recall of data from the more artistic charts). Note they define "**ratio**" different from how we have defined in class.