# Final Project and Proposal

https://iu.instructure.com/courses/22166
27/assignments/15873434

# Data Visualization

for interactive (e.g., Plotly) **and** static (e.g., Seaborn) visualizations

# What is Data Visualization

*"The use of computer-generated, interactive, **visual representations of data** to amplify cognition."*

*"The **transformation of data into visual representations** to aid people in the analysis, exploration, and communication of that data"*

Card et al. 1999 and Jeff Heer

# Data

What are we gathering?

# Types of Data

- Healthcare

- Financial Markets

- Scientific Data

- Social Media

- etc.

# How much data are we gathering

- Data on Wikipedia *alone* in 2023 = **439.17 TB**

    - Over 100 million files.

- This is just one (relatively lightweight) server

- Extrapolate this to the whole of the internet …

- 1000 "Wikipedias" would take up 420 petabytes

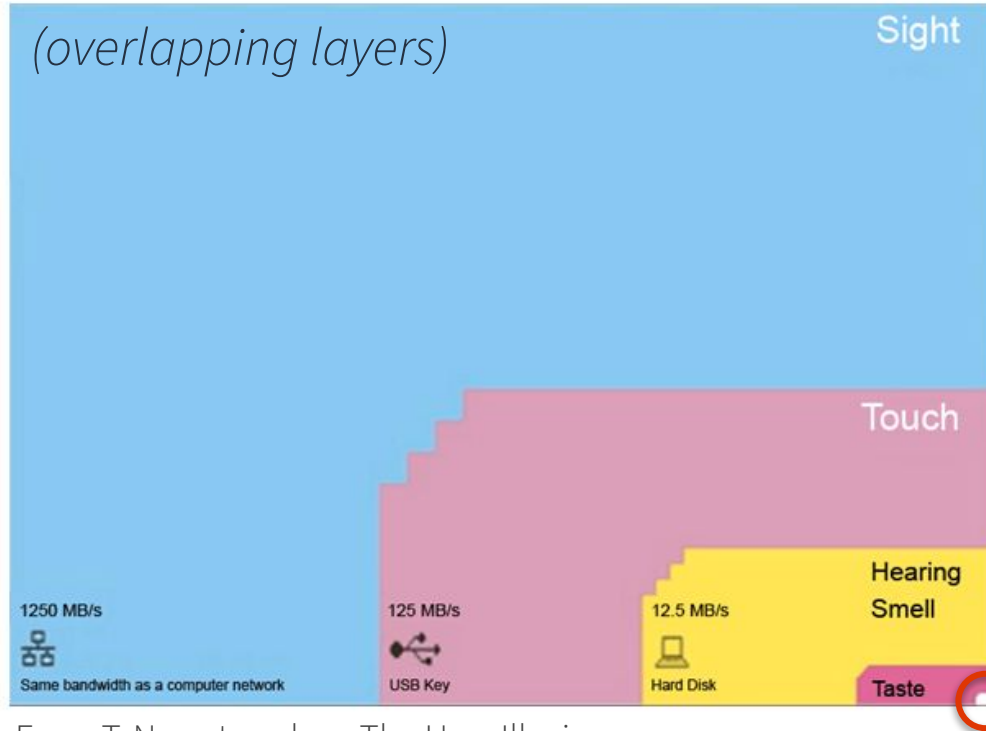    - 420,… with 16 zeros … bytes (e.g., an integer takes up ~4 bytes)

# How to interpret it?

Sight, Hearing, Touch, Smell, Taste, Telepathy

# How to interpret it?



**Sight**,
Hearing,
Touch,
Smell,
Taste,
Telepathy

# Sensory Bandwidth



*(overlapping layers)*

Sight

Touch

Hearing
Smell

1250 MB/s
Same bandwidth as a computer network

125 MB/s
USB Key

12.5 MB/s
Hard Disk

Taste

This small white box is what we're actually conscious of …

From T. Norretranders, The User Illusion:
*Cutting Consciousness Down to Size*, 1999

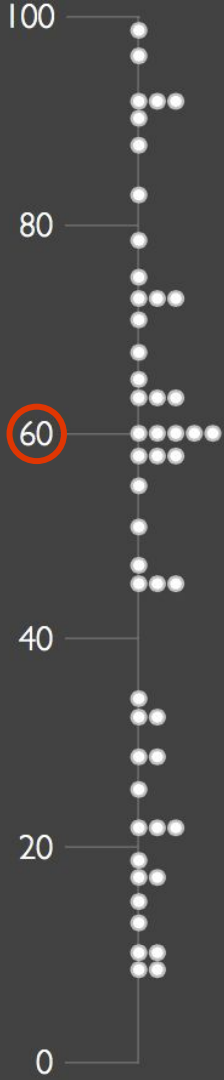# How many "X" letters are here?

345OIJDFG98C90U5ET09VBKK23490XIVBCIBJ0345T09U
2G84GDF09U34590IDFK90345I-09345K90FU90DF90JDF
34T09X90DFJG90J34T09J34509J3459DFG08JKLSTJP435
DFDFG45OJERPOTJ45OPIJFDGLKM34T5XJSCTYY7K456
POJ345OIJLGJKOPE390UVFHUDGH9345H9R4N97HWTIO
MADSIOPEJDFGPJ4309UT509345PODFGX093490823JFD
PWDEIJ3408UDFMV984385Y0834N92384YU8DFB0H3T4N
345J09JDFG09J345X98U5Y09JGFB089H34509UJ45TM0IG
P5JDGIOEGWJPIO345U345OPIJDTOPI3458345JPODFG09
45POJ34X09345J08EFJ825HJDFSJIPADOPQWIXERWNVF

345OIJDFG98C90U5ET09VBKK23490**X**IVBCIBJ0345T09U
2G84GDF09U34590IDFK90345I-09345K90FU90DF90JDF
34T09**X**90DFJG90J34T09J34509J3459DFG08JKLSTJP435
DFDFG45OJERPOTJ45OPIJFDGLKM34T5**X**JSCTYY7K456
POJ345OIJLGJKOPE390UVFHUDGH9345H9R4N97HWTIO
MADSIOPEJDFGPJ4309UT509345PODFG**X**093490823JFD
PWDEIJ3408UDFMV984385Y0834N92384YU8DFB0H3T4N
345J09JDFG09J345**X**98U5Y09JGFB089H34509UJ45TM0IG
P5JDGIOEGWJPIO345U345OPIJDTOPI3458345JPODFG09
45POJ34**X**09345J08EFJ825HJDFSJIPADOPQWI**X**ERWNVF

| 15 | 19 | 60 |
|----|----|----|
| 33 | 11 | 75 |
| 57 | 34 | 79 |
| 18 | 51 | 92 |
| 73 | 22 | 13 |
| 71 | 60 | 22 |
| 17 | 10 | 68 |
| 73 | 18 | 55 |
| 65 | 46 | 29 |
| 60 | 73 | 22 |
| 46 | 92 | 97 |
| 10 | 58 | 46 |
| 57 | 17 | 83 |
| 26 | 99 | 33 |
| 88 | 92 | 60 |
| 91 | 29 | 57 |
| 96 | 12 | 47 |

Given these 50 numbers
what number appears most often?

Given these 50 numbers
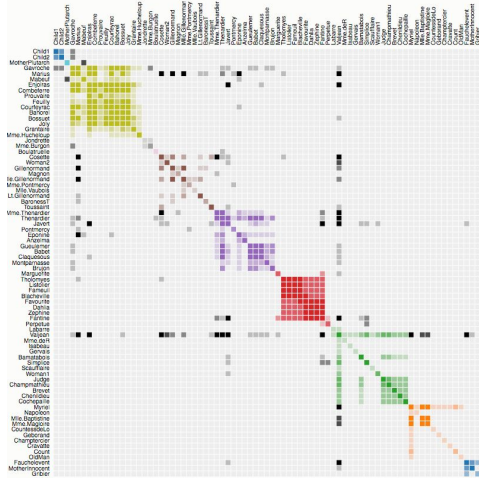what number appears most often?

# Visuals in our Culture

- "I see what you're saying"

- "Seeing is believing"

- "I now see the big picture"

- "A picture is worth a thousand words"

# Computers vs. People

- Computers

    - Process large quantities of information quickly

    - Computing a specific task

    - Re-use of code for different datasets

- Humans

    - Devising questions when we don't know what to look for

    - We can use visuals to help with this

# Building Visualizations

# Spectrum of Visualization Use



Profits Doubled on Treated Group within 6 Months

**Analysis**
*Better understand data*

**Communication**
*tell a story with data*

# Spectrum of Visualization Use



**Analysis**

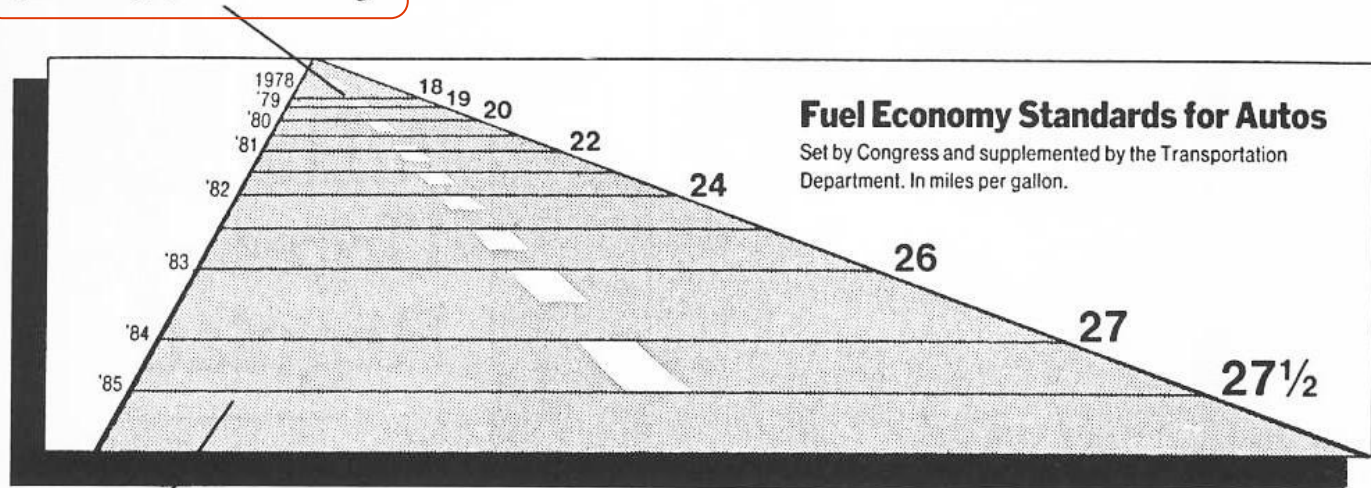*ad hoc; for your eyes only*

**Communication**

*more presentable!*

# Visual Integrity: **not to lie with data**

The representation of numbers, as physically measured on the surface of the graphic itself, should be directly proportional to the quantities represented.

This line, representing 18 miles per gallon in 1978, is 0.6 inches long.

This line, representing 27.5 miles per gallon in 1985, is 5.3 inches long.

**Fuel Economy Standards for Autos**

Set by Congress and supplemented by the Transportation Department. In miles per gallon.

$$\text{lie factor} = \frac{\text{graphic effect}}{\text{numeric effect}}$$

$$= \frac{(5.3 - 0.6)/0.6}{(27.5 - 18)/18} = \frac{7.8}{0.5} = 15.6$$

# Plot Types

# Big Numbers

- Great to display simple information

- Tips:

  - Include a baseline for comparison (e.g., an average)

Page Views:

# 5,567

Daily Average:  3,625

# Tables

- Conveys comparisons across categories

- Tips:
  - Too big is less effective
  - Use background color sparingly
  - Bold what's important (e.g., the data)
  - Keep to less than 3 x 3
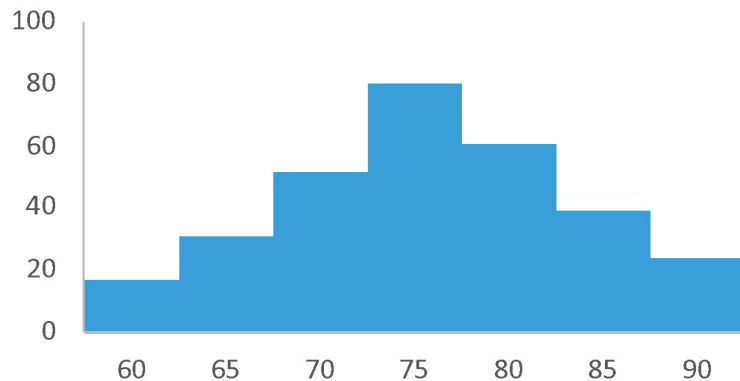
|  | Likes Chocolate | Dislikes Chocolate |
|---|---|---|
| Children | 65 | 5 |
| Adults | 40 | 30 |

# Bar Charts

- Compare different groups or categories

- Tips:

  - Sort based on something meaningful (e.g. length, alphabetical)

  - **Start scale at 0 (for all plots, really)**

  - Horizontal, vertical and stacked bars are commonly used

# Histograms

- Represents the distribution of data

- Tips:

  - Range of values must be binned

  - Can be normalized
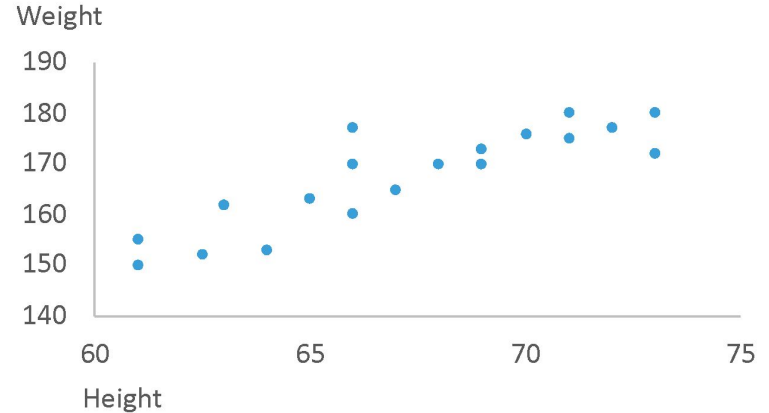
    (sum of bar heights is 1)

# Line Charts

- Track changes over time

- Height and slope lets us see trends

- Tips:

  - x-axis should be *continuous* data

  - Time is represented from left to right
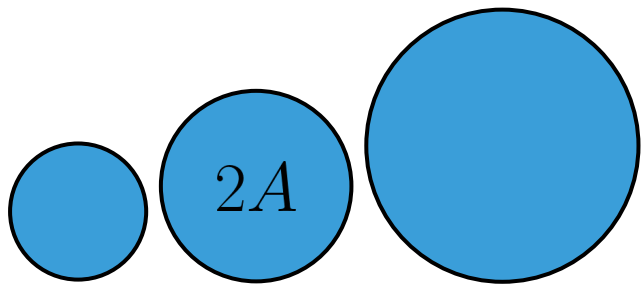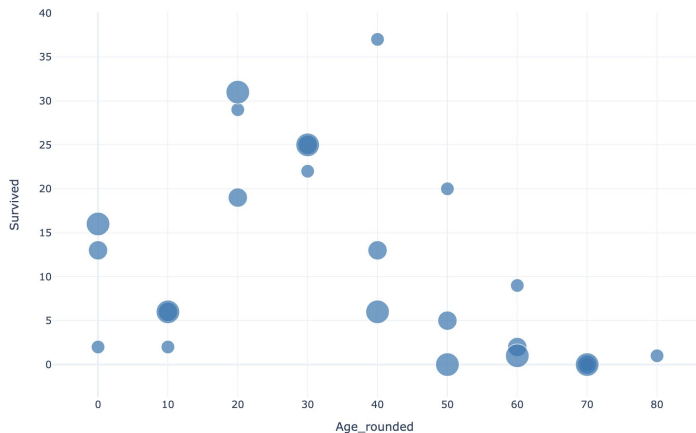
  - **Do not use if x-axis is not ordered**

Price ($)

# Scatter Plots

- Displays relationship between two measures

- Tips:

  - Best for continuous data

  - Avoid using with qualitative data

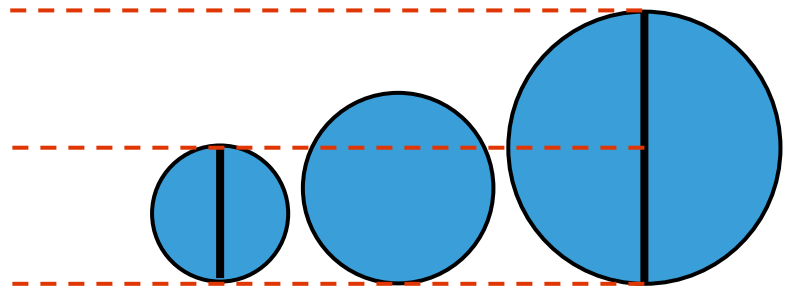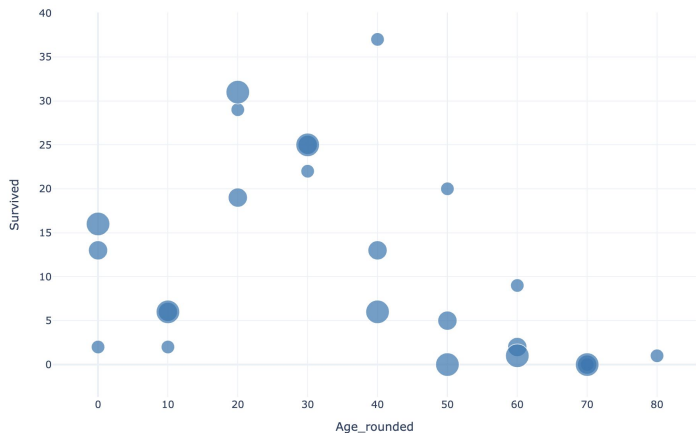# Bubble Charts

- Used for scatter plots with 3 variables

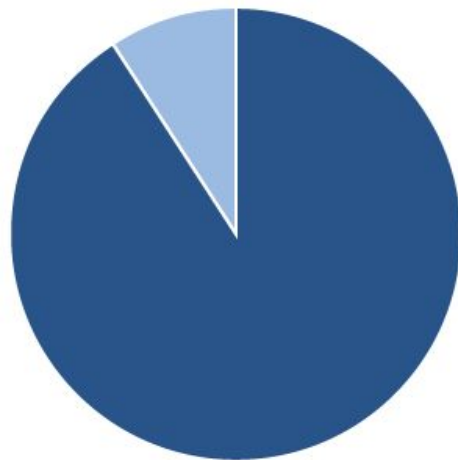- Tips:

  - Area causes confusion

  - Use diameter

# Bubble Charts

- Used for scatter plots with 3 variables

- Tips:

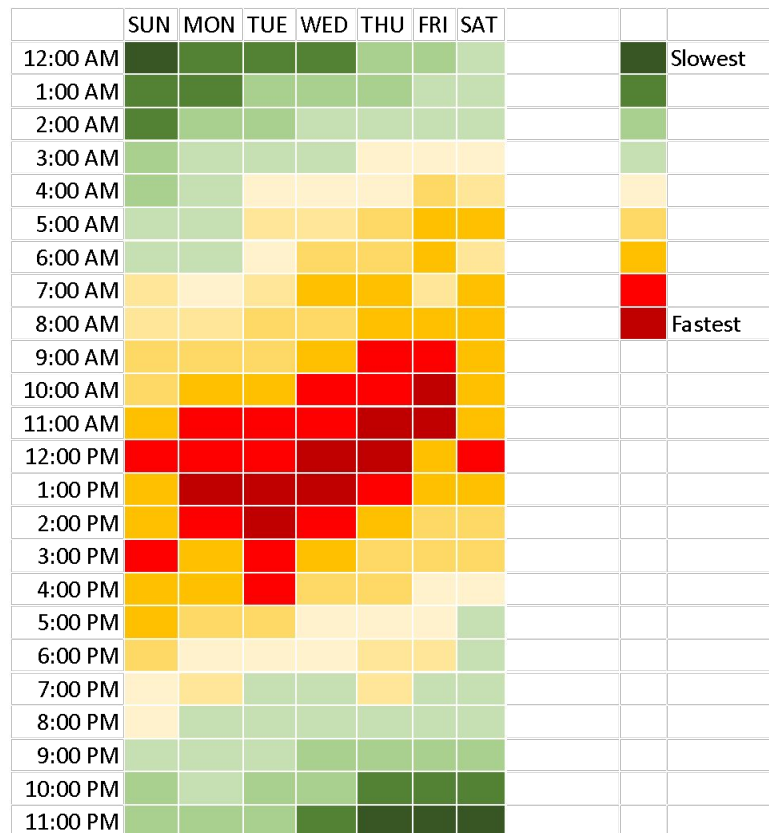  - Area causes confusion

  - Use diameter

# Pie Charts (think "slice")

- Compare parts of a whole

- Tips:

  - All parts must sum to 100%

  - Best for binary data

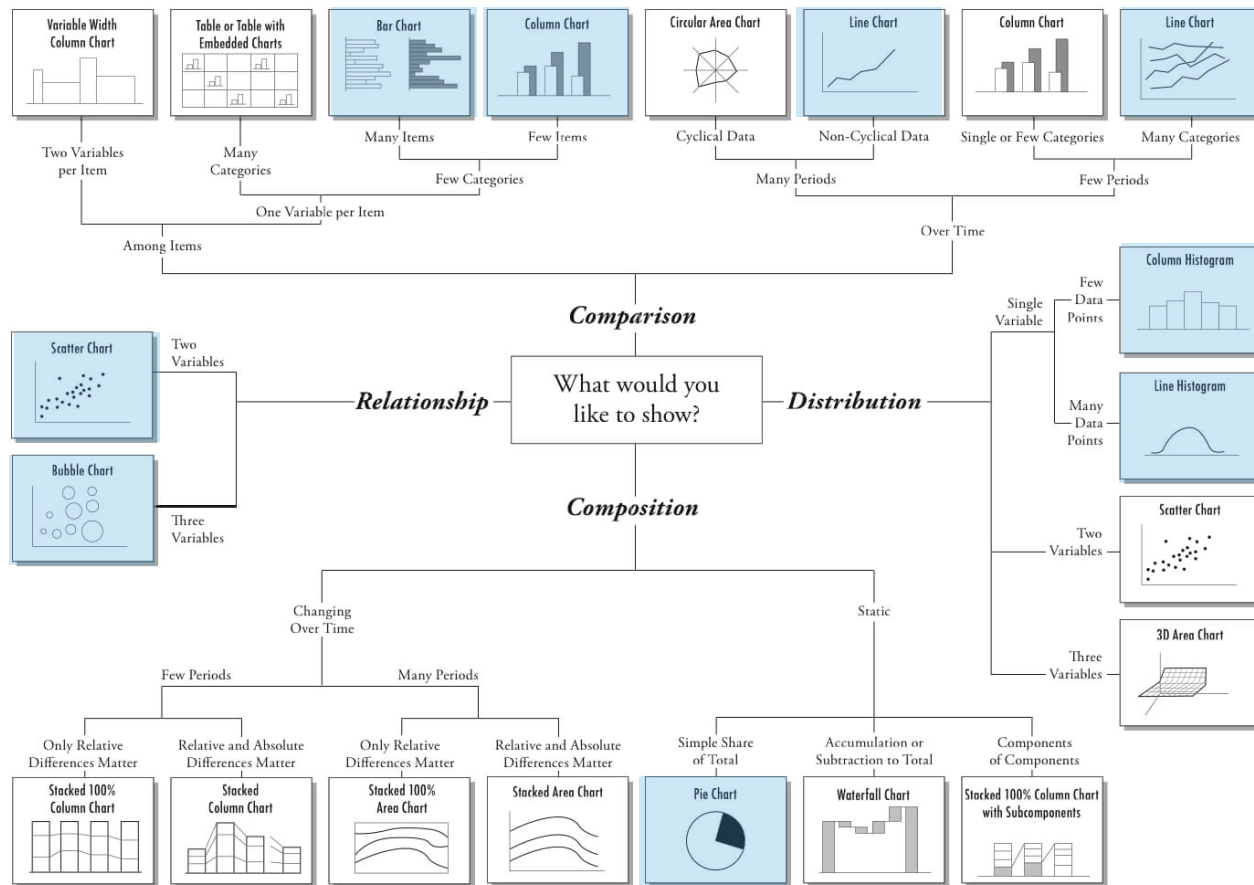- **Only use if one is very small**

# Heat Maps

- Use color to visualize a matrix

- Tip:

  - **Best when x *and* y variables are ordinal** (E.g., weekday, low to high, distance)

  - Use sensible color maps!

  - Color *by-column* or *by-row*

- Very useful for binary data (e.g., missing)

|          | SUN | MON | TUE | WED | THU | FRI | SAT |     |     |          |
|----------|-----|-----|-----|-----|-----|-----|-----|-----|-----|----------|
| 12:00 AM |     |     |     |     |     |     |     |     |     | Slowest  |
| 1:00 AM  |     |     |     |     |     |     |     |     |     |          |
| 2:00 AM  |     |     |     |     |     |     |     |     |     |          |
| 3:00 AM  |     |     |     |     |     |     |     |     |     |          |
| 4:00 AM  |     |     |     |     |     |     |     |     |     |          |
| 5:00 AM  |     |     |     |     |     |     |     |     |     |          |
| 6:00 AM  |     |     |     |     |     |     |     |     |     |          |
| 7:00 AM  |     |     |     |     |     |     |     |     |     |          |
| 8:00 AM  |     |     |     |     |     |     |     |     |     | Fastest  |
| 9:00 AM  |     |     |     |     |     |     |     |     |     |          |
| 10:00 AM |     |     |     |     |     |     |     |     |     |          |
| 11:00 AM |     |     |     |     |     |     |     |     |     |          |
| 12:00 PM |     |     |     |     |     |     |     |     |     |          |
| 1:00 PM  |     |     |     |     |     |     |     |     |     |          |
| 2:00 PM  |     |     |     |     |     |     |     |     |     |          |
| 3:00 PM  |     |     |     |     |     |     |     |     |     |          |
| 4:00 PM  |     |     |     |     |     |     |     |     |     |          |
| 5:00 PM  |     |     |     |     |     |     |     |     |     |          |
| 6:00 PM  |     |     |     |     |     |     |     |     |     |          |
| 7:00 PM  |     |     |     |     |     |     |     |     |     |          |
| 8:00 PM  |     |     |     |     |     |     |     |     |     |          |
| 9:00 PM  |     |     |     |     |     |     |     |     |     |          |
| 10:00 PM |     |     |     |     |     |     |     |     |     |          |
| 11:00 PM |     |     |     |     |     |     |     |     |     |          |

# Design

# Chart Suggestions—A Thought-Starter

Variable Width Column Chart — Two Variables per Item

Table or Table with Embedded Charts — Many Categories

Bar Chart — Many Items

Column Chart — Few Items

Few Categories

One Variable per Item

Among Items

Circular Area Chart — Cyclical Data

Line Chart — Non-Cyclical Data

Many Periods

Column Chart — Single or Few Categories

Line Chart — Many Categories

Few Periods

Over Time

**Comparison**

**Relationship**

Scatter Chart — Two Variables

Bubble Chart — Three Variables

## What would you like to show?

**Distribution**

Column Histogram — Few Data Points

Line Histogram — Many Data Points

Single Variable

Scatter Chart — Two Variables

3D Area Chart — Three Variables

**Composition**

Changing Over Time

Few Periods

Only Relative Differences Matter — Stacked 100% Column Chart

Relative and Absolute Differences Matter — Stacked Column Chart

Many Periods

Only Relative Differences Matter — Stacked 100% Area Chart

Relative and Absolute Differences Matter — Stacked Area Chart

Static

Simple Share of Total — Pie Chart

Accumulation or Subtraction to Total — Waterfall Chart

Components of Components — Stacked 100% Column Chart with Subcomponents

https://extremepresentation.typepad.com/files/choosing-a-good-chart-09.pdf

# Visual Principles

# Text Orientation

Vertical text orientation requires some mental effort!

# Text Orientation

Vertical text orientation requires mental effort!

So does this!

# Text Orientation

Vertical text orientation requires mental effort!

So does this!

Horizontal text is easier to read!

# Font Size

# Font Size

# Labels

**Always** include all descriptive labels:

- Title of Plot
- Axis labels
- Legend (color, line, etc.)
- Highlights and Callouts



Population in a Fake Country During a Made-up Time Period

**Leave nothing ambiguous or unclear!**

# Color

# Color

Sequential

Diverging

Categorical

Highlight

# Color

Sequential  Pick **one** hue; captures "amount"

Diverging  Middle point *must* hold **meaning** (e.g., 0)

Categorical  No color should stand out

Highlight  Communicate w/ color (e.g., red = bad)

# Color Palette

- Limit categorical colors to 5

  - 7 at the absolute most.

# Color consistency across charts



No color should stand out more than any other.

# RGB to grayscale

Grayscale conversion "equates" some hues …

RGB Image

Grayscale Image

# RGB to grayscale

Grayscale conversion "equates" some hues …

# Color blindness

0.5% 8.0%

Instead of red and green
use blue and orange.

Consider color-blind color scale
(e.g., "viridis")

# Count the 2s

5498731840

4893128612

2634085106

1592059852

3854634876

# Count the 3s

549873<span style="color:red">3</span>1840

489<span style="color:red">3</span>128612

26<span style="color:red">3</span>4085106

1592059852

<span style="color:red">3</span>85463<span style="color:red">3</span>4876

# Balance background and foreground

ACVLSIGBSLWUHKAJSLHV

ACVLSIGBSLWUHKAJSLHV

ACVLSIGBSLWUHKAJSLHV

# Pre-attentive attributes



Orientation

Width

Size

Shape

Color

Enclosure

We perceive these differences unconsciously before we are aware of them

# What is the main point?

# What is the main point?

Highlight your conclusion

Non-critical information should be removed or put in the background

# Gridlines



Keep when illustrating specific numeric values

Remove when the "message" is the trend itself

# Increase Data-Ink Ratio

Aim for a high Data-to-Ink Ratio:
- Maximize visual elements related explicitly to data
- Minimize extra pixels not related to data (e.g., background color, grids)

# 30 second rule

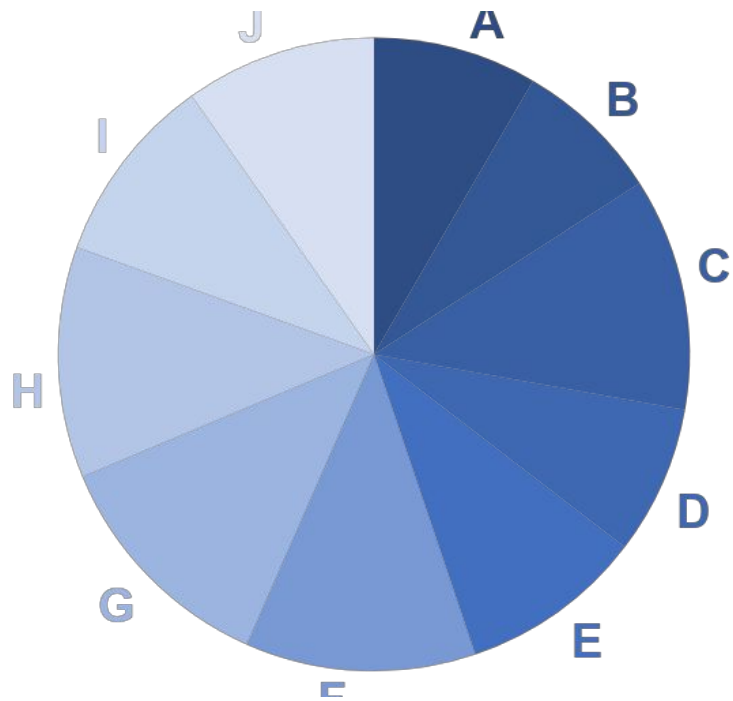A viewer should be able to interpret the message of a visualization within 30 seconds

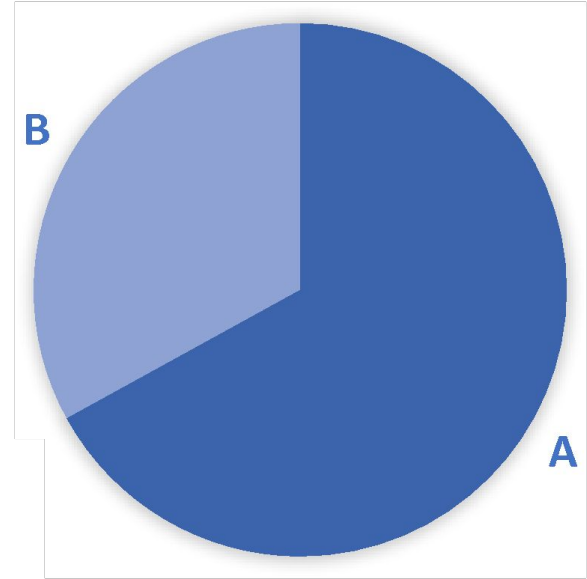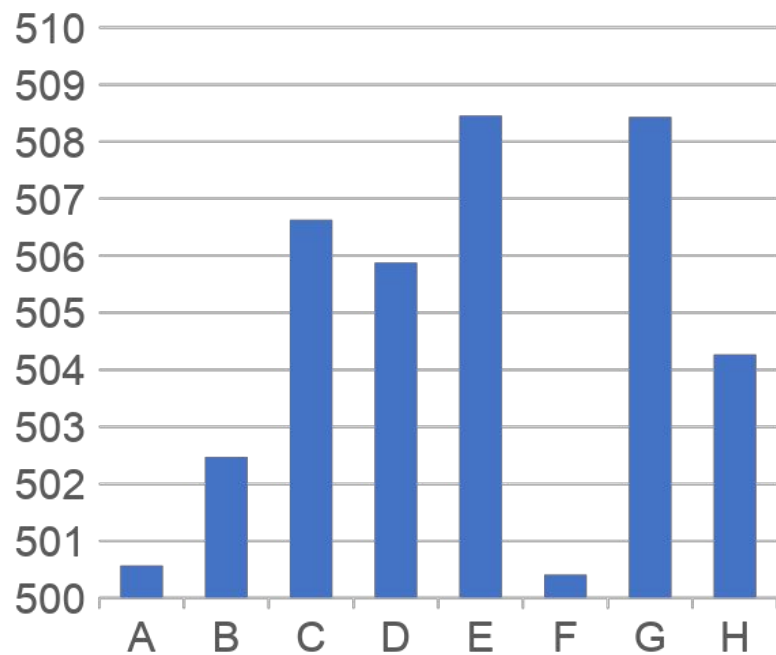# Common Mistakes

# Avoid 3D

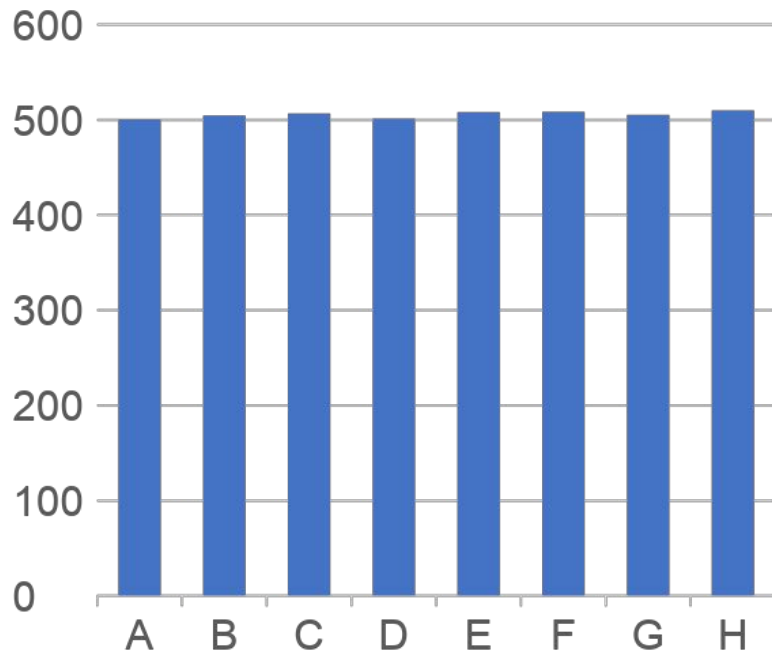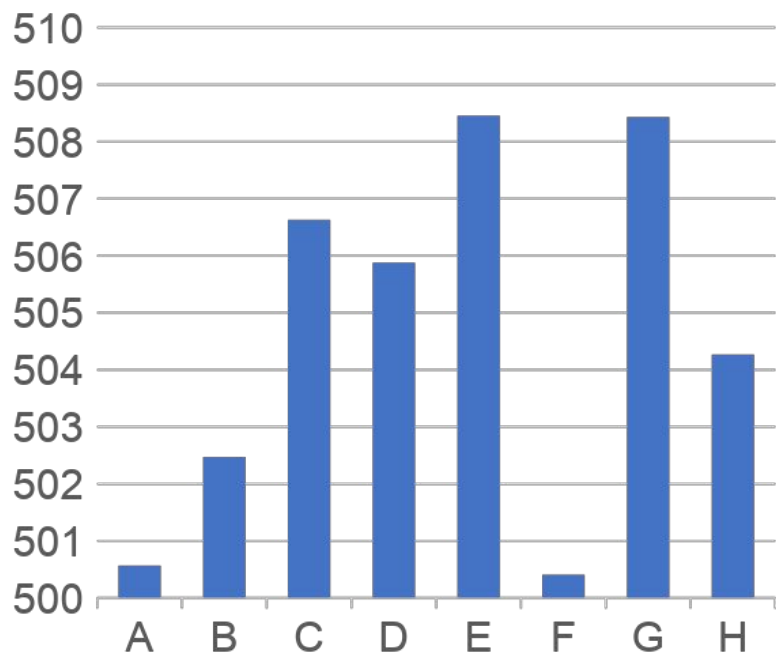# Which is the biggest slice?

# Sometimes pie charts work well

Keep to 2 slices, maximum

# "Scale" requires a zero value
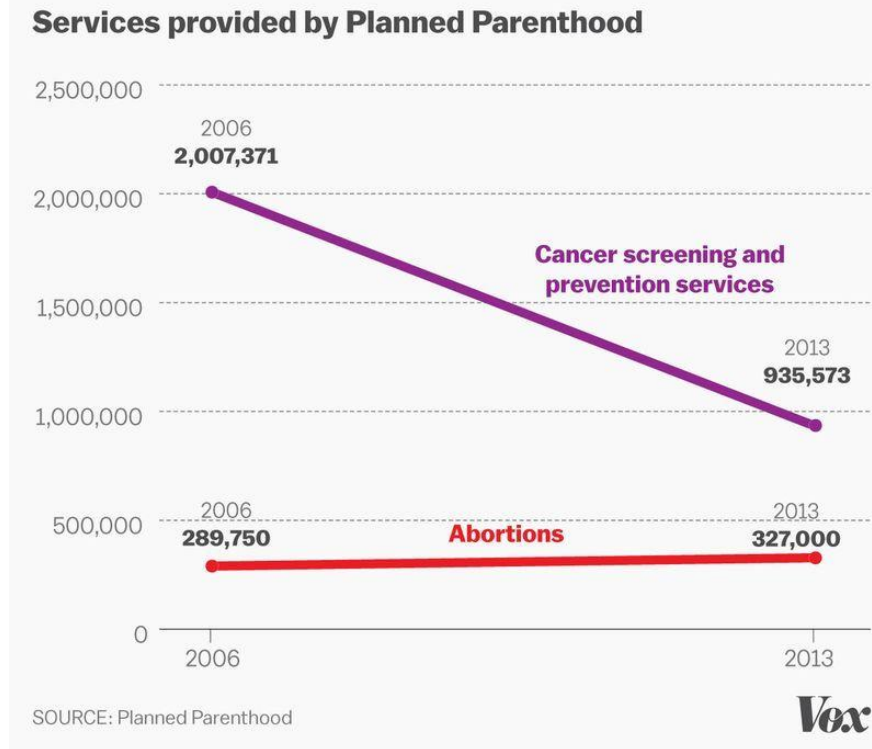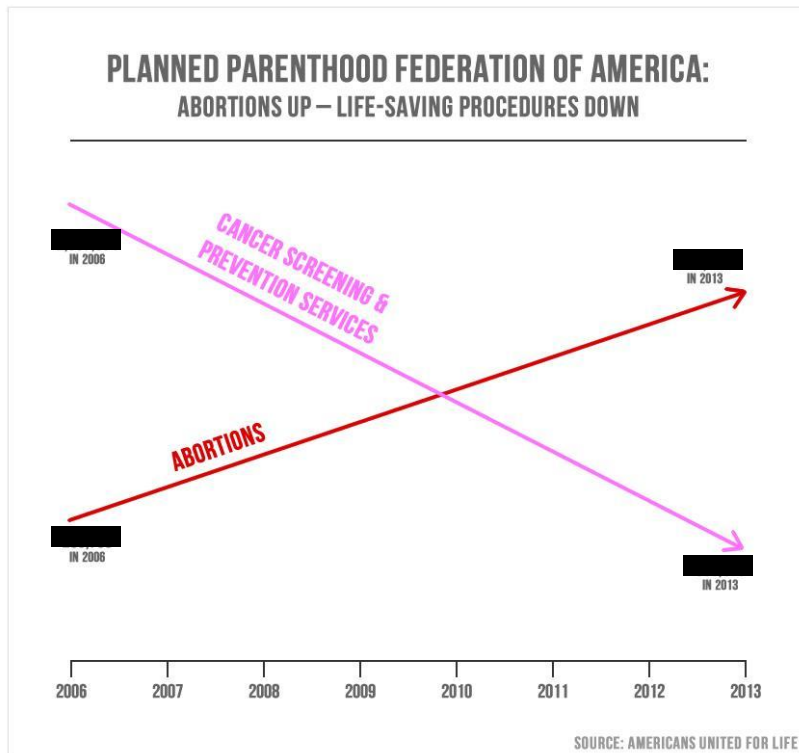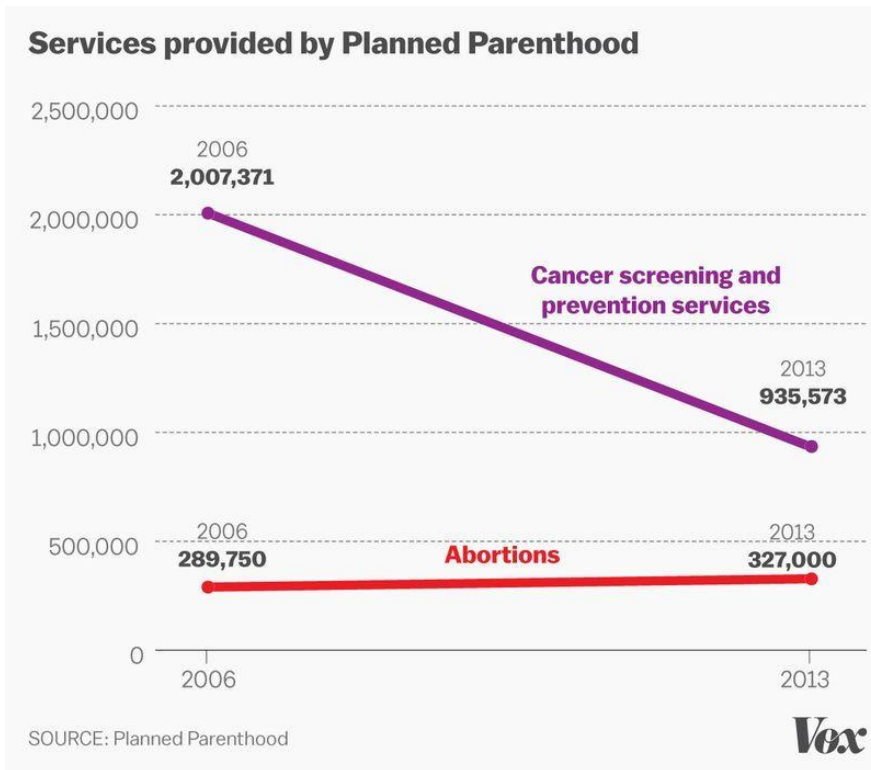
# "Scale" requires a zero value



**Start numerical axes at zero (or include a zero line)**

# Example



**PLANNED PARENTHOOD FEDERATION OF AMERICA:**
ABORTIONS UP — LIFE-SAVING PROCEDURES DOWN

CANCER SCREENING & PREVENTION SERVICES

ABORTIONS

IN 2006   IN 2013   IN 2006   IN 2013

2006   2007   2008   2009   2010   2011   2012   2013

SOURCE: AMERICANS UNITED FOR LIFE

**Services provided by Planned Parenthood**

2,500,000

2006
**2,007,371**

2,000,000

**Cancer screening and prevention services**

1,500,000

2013
**935,573**

1,000,000

500,000

2006
**289,750**

**Abortions**

2013
**327,000**

0

2006   2013

SOURCE: Planned Parenthood

*Vox*

# Example

Synchronize the scale of **both** axes



**Services provided by Planned Parenthood**

2,500,000

2006
**2,007,371**

2,000,000

**Cancer screening and prevention services**

1,500,000

2013
**935,573**

1,000,000

2006
**289,750**

**Abortions**

2013
**327,000**

500,000

0

2006

2013

SOURCE: Planned Parenthood

*Vox*

# Example

Avoid inverted axes

# Bad Examples

# What is wrong with this visual?

- Readability

- Colors



**Types of debt**

The total owed by the average U.S. household, by debt type.

| | |
|---|---|
| Credit cards | $16,748 |
| Mortgages | $176,222 |
| Auto loans | $28,948 |
| Student loans | $49,905 |
| Any type of debt | $134,643 |

# What is wrong with this visual?

- 3D

- Too many slices

- >100%



Die am häufigsten genannten bewährten Gegenstände an Bord

- Generator 8%
- Log/Echolot 8%
- Windsteuerung 52%
- Kühlschrank 9%
- SSB Radio 35%
- Satellitentelefon 11%
- Solarzellen 31%
- Rollanlage (Vorsegel) 12%
- Autopilot 29%
- Außenborder 12%
- Wassermacher 25%
- Ankerwinde 13%
- Dieselmotor 20%
- Radar 14%
- Windgenerator 20%
- Bügelanker 14%
- Dingi 19%
- AIS 14%
- Segel 19%
- GPS 15%

# What is wrong with this visual?

- Decimals

- y-axis



**LEARNING HOW TO LEARN**
Percentage of VC-backed companies with certain words in their company description over time

artificial intelligence — deep learning — machine learning

CBINSIGHTS

# What is wrong with this visual?

- 30-second rule

- … What is happening?



The ₿bitcoin Wealth Distribution

4.11% OF ADDRESSES OWN 96.53% OF BTC*

0.00088% of the addresses own 17.49% of BTC

0.01% of the addresses own 20.47% of BTC

0.94% of the addresses own 28.02% of BTC

0.10% of the addresses own 21.90% of BTC

3.06% of the addresses own 7.92% of BTC

0.00000748% of the addresses own 0.73% of BTC

41.93 % of the addresses own 0.01% of BTC

24.94% of the addresses own 0.09% of BTC

9.41% of the addresses own 2.84% of BTC

19.61% of the addresses own 0.53% of BTC

95.89% OF ADDRESSES OWN 3.47% OF BTC*

* Data as of September 12th, 2017
Article and Sources:
https://howmuch.net/articles/bitcoin-wealth-distribution
https://bitcoinprivacy.net/

howmuch .net

# Oh no …

# Tools

# Visualization Tools

|              | Excel   | Seaborn | Plotly      | Tableau     | D3              |
|--------------|---------|---------|-------------|-------------|-----------------|
| Interactivity | Static  | Static  | Interactive | Interactive | Interactive     |
| Difficulty   | Easy    | Medium  | Medium      | Easy        | Difficult       |
| Pros         | Popular | Python  | Python      | Beautiful   | Custom          |
| Cons         | Basic   | Python  | Python      | $$$         | Time Consuming  |

# Visualization Tools

| | Excel | Seaborn | Plotly | Tableau | D3 |
|---|---|---|---|---|---|
| Interactivity | Static | Static | Interactive | Interactive | Interactive |
| Difficulty | Easy | Medium | Medium | Easy | Difficult |
| Pros | Popular | Python | Python | Beautiful | Custom |
| Cons | Basic | Python | Python | $$$ | Time Consuming |

# Visualization Makeover

# Example

Use line chart!



Profit per Group by Month

# Example

Remove background color!

# Example

Start axis at 0!



Profit per Group by Month

# Example



Decimal points!

# Example

Vertical "Profit" label!

# Example

Diagonal x-axis!
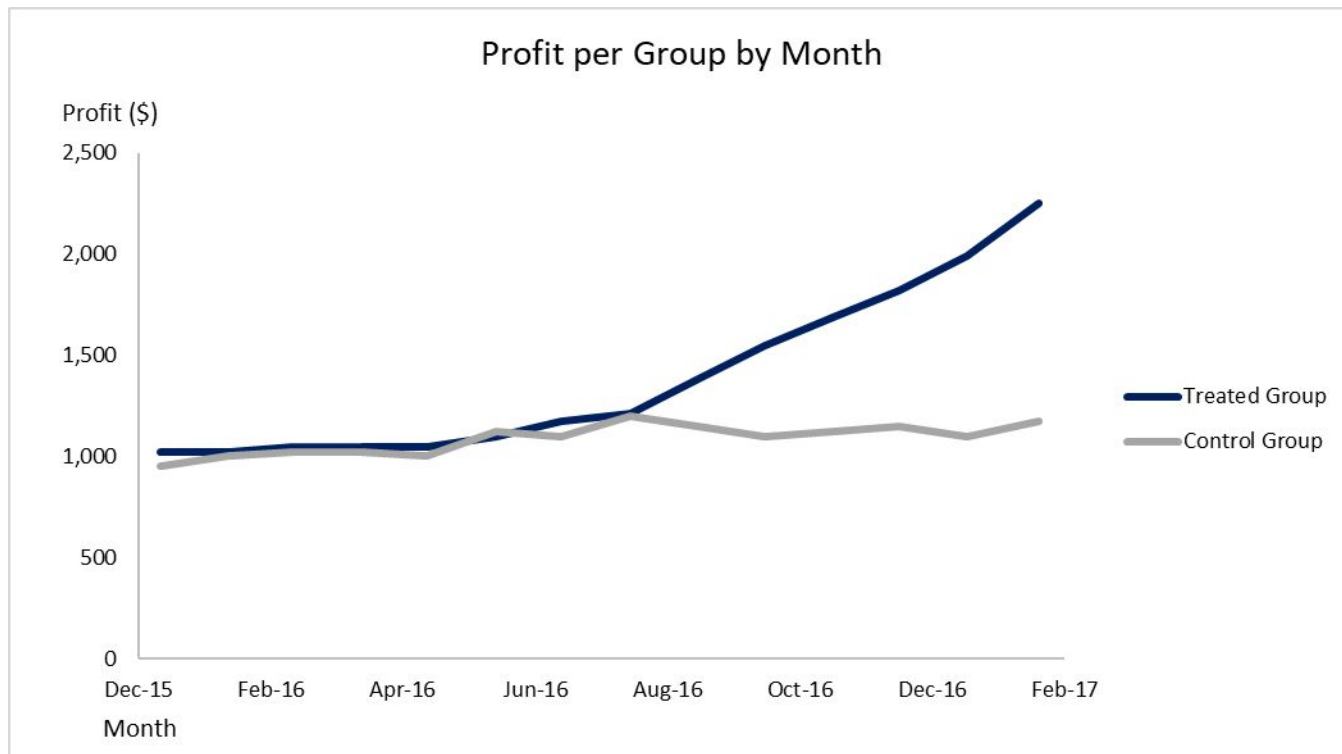
# Example

X-Label left justified!

# Example

Gridlines!

# Example

Color and dots! (this is a *binary* time series)

# Example

Legend, proximity and space



Profit per Group by Month

# Example

Axis background



Profit per Group by Month

# Example

Descriptive title and left!

# Example

Profits Doubled on Treated Group within 6 Months
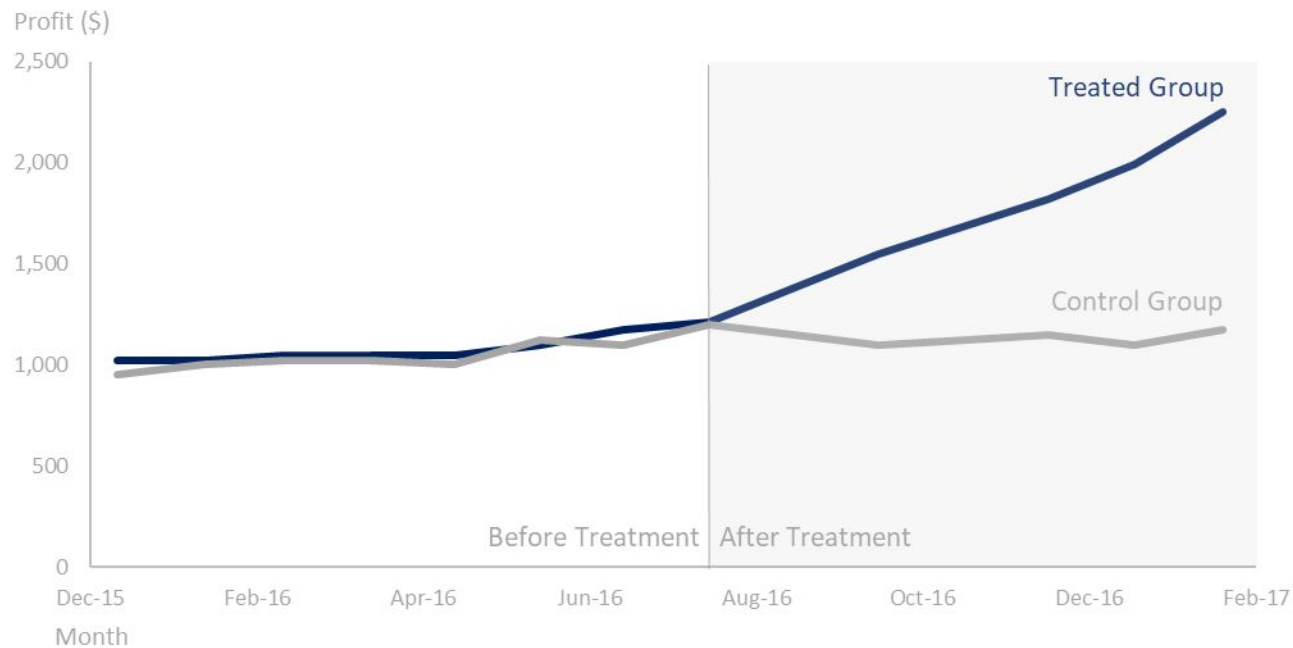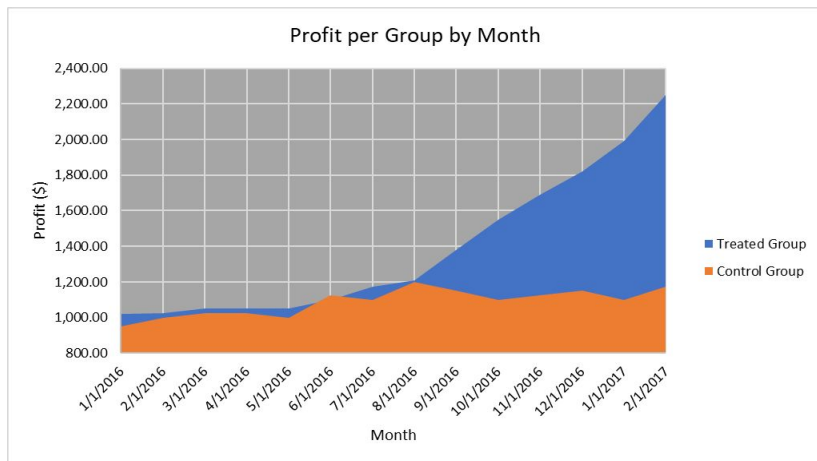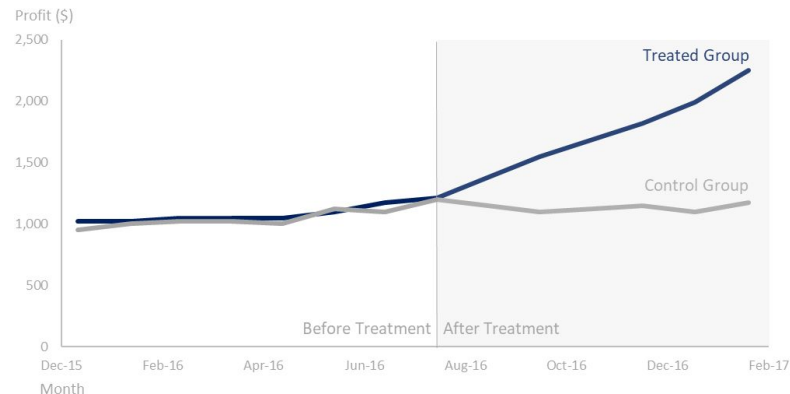
# Example

Box!



Profits Doubled on Treated Group within 6 Months

# Example



Profits Doubled on Treated Group within 6 Months

# Example

# Summary

- Keep design thinking in mind

  - Charts should be easy to understand

  - Visuals should be discoverable

- Use text, color, and highlighting techniques

- Avoid unnecessary ink

  - Avoid 3D

  - Pie charts with too many slices

- Don't lie with data

- Don't use defaults – spend some time making over your visualizations

*"Above all else, show the data."*

~ Tufte, 1983