

# Graphing Best Practices

CSS 2 – Spring, 2022

Erik Brockbank

\*content heavily borrowed from a graduate lecture by Edward Vul

# What We'll Cover

# What We'll Cover

1. Ways graphs can go wrong (*there are more than you think...*)

# What We'll Cover

1. Ways graphs can go wrong (*there are more than you think...*)
2. How to avoid these pitfalls

# What We'll Cover

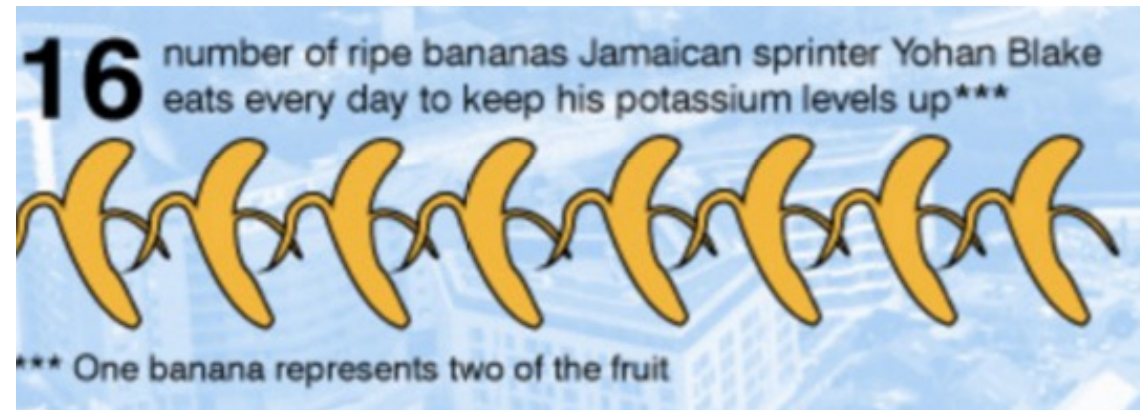
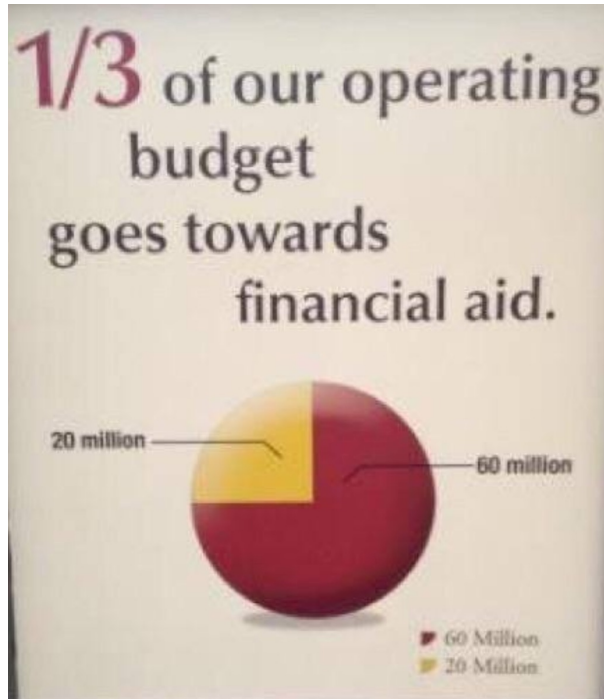
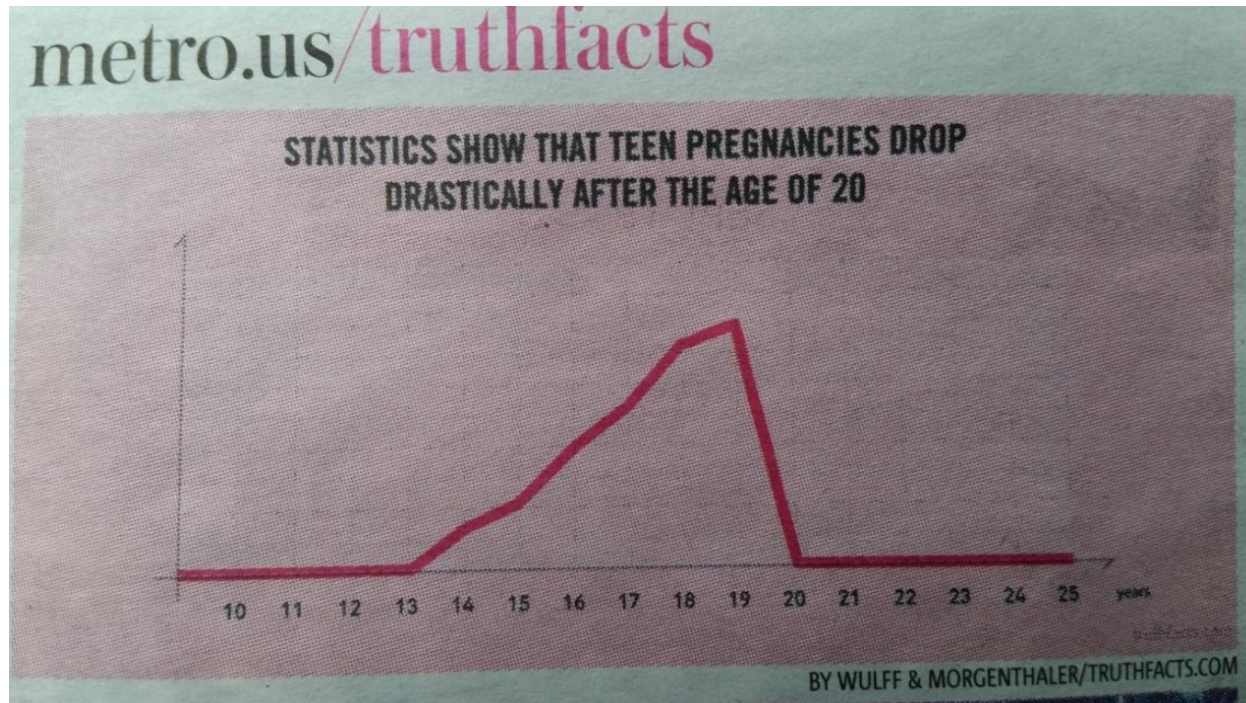
1. Ways graphs can go wrong (*there are more than you think...*)
2. How to avoid these pitfalls
3. Choosing the right visualization

# What We'll Cover

1. Ways graphs can go wrong (*there are more than you think...*)
2. How to avoid these pitfalls
3. Choosing the right visualization

# Graphs Gone Awry

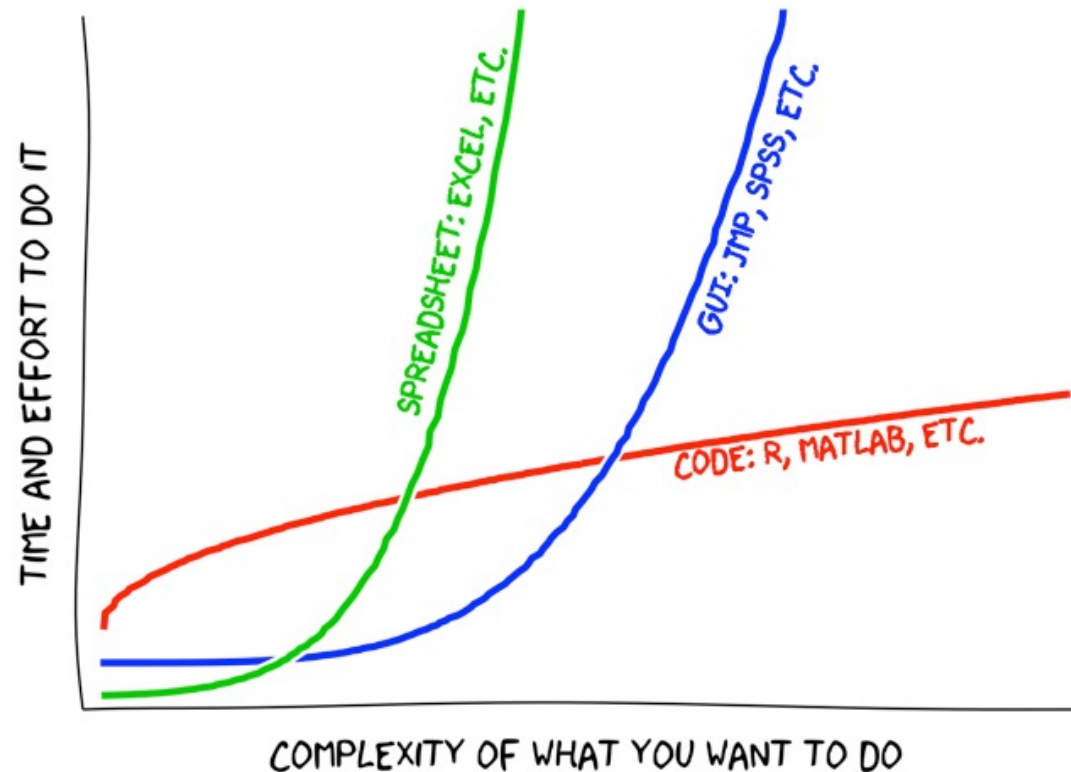
# Graphs Gone Awry





# Graphs Gone Awry

Sometimes things that look like graphs are made up...



# Graphs Gone Awry

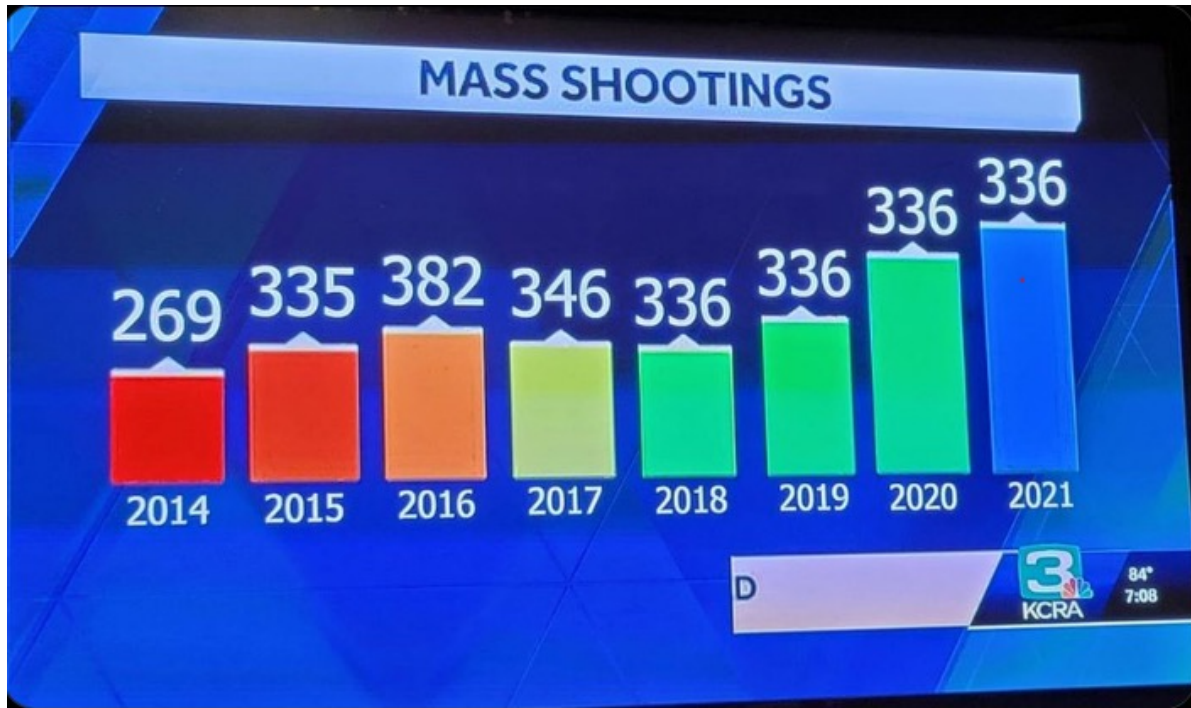
Sometimes things that look like graphs are made up...



source: @GraphCrimes

# Graphs Gone Awry

Sometimes it's a real graph but it doesn't match the data...



# Graphs Gone Awry

Axis scales can be misleading...

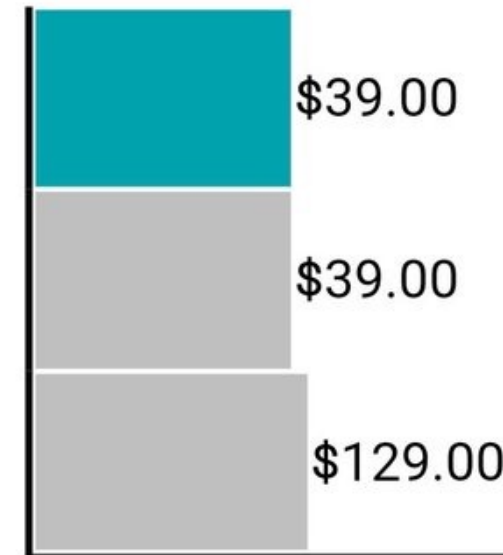


This bill compared to previous months

17 Oct

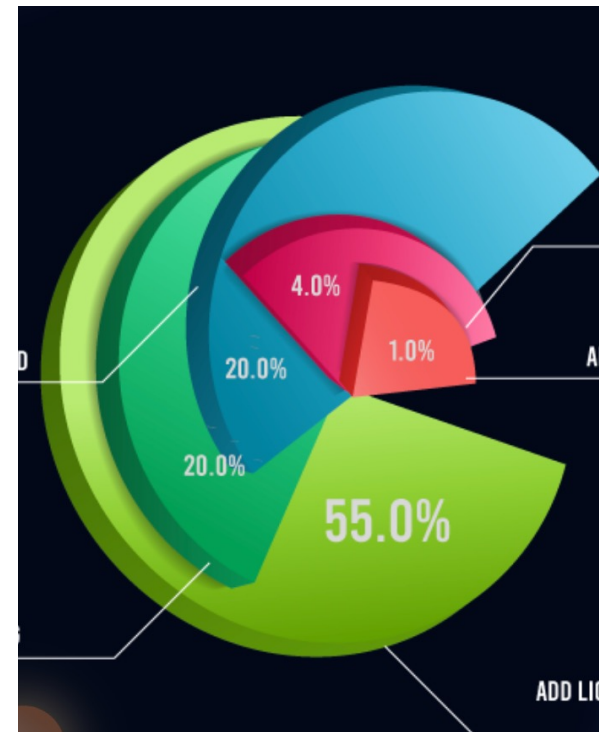
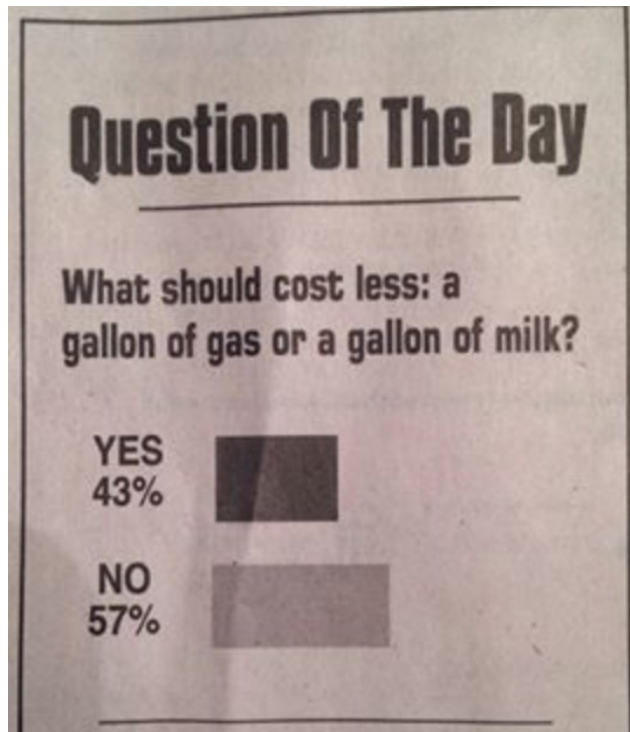
17 Sep

17 Aug



# Graphs Gone Awry

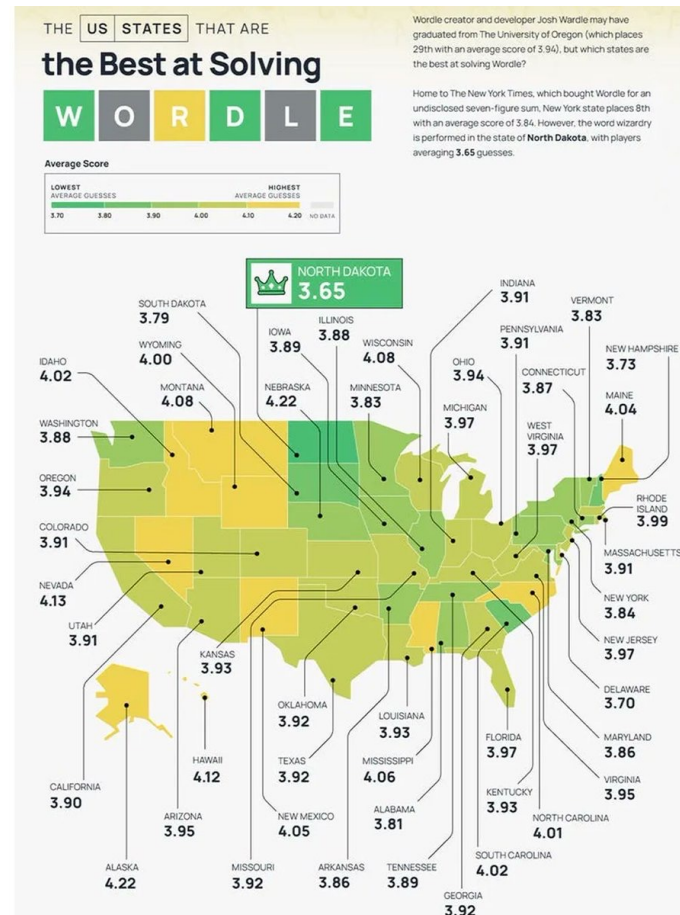
...or not there at all





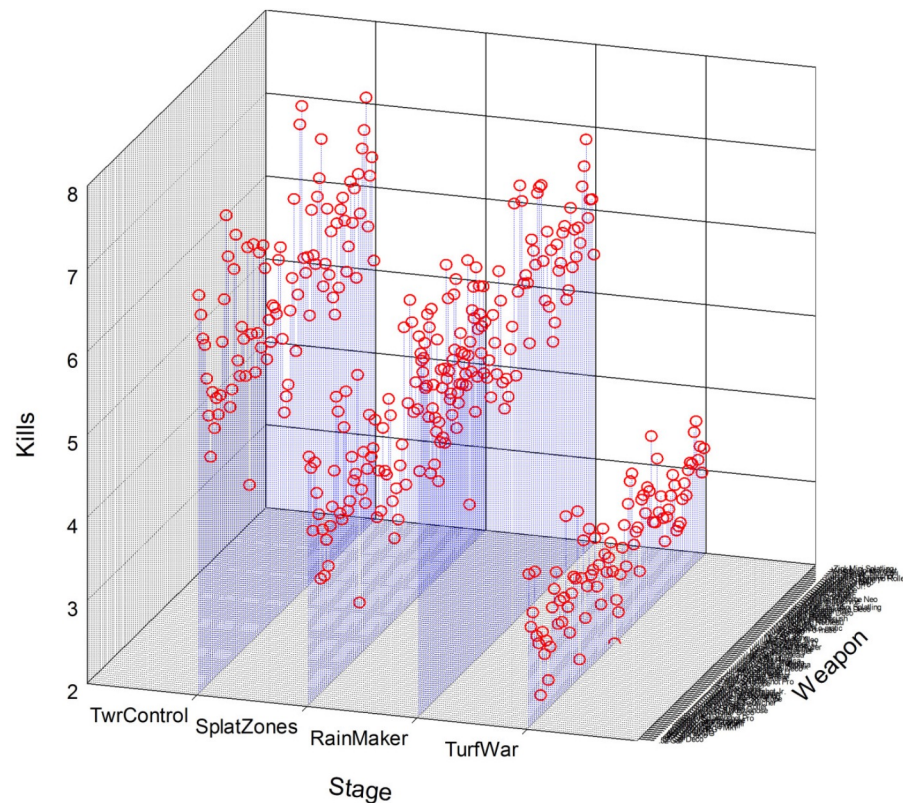
# Graphs Gone Awry

Be careful of bins or gradients that are kind of useless...

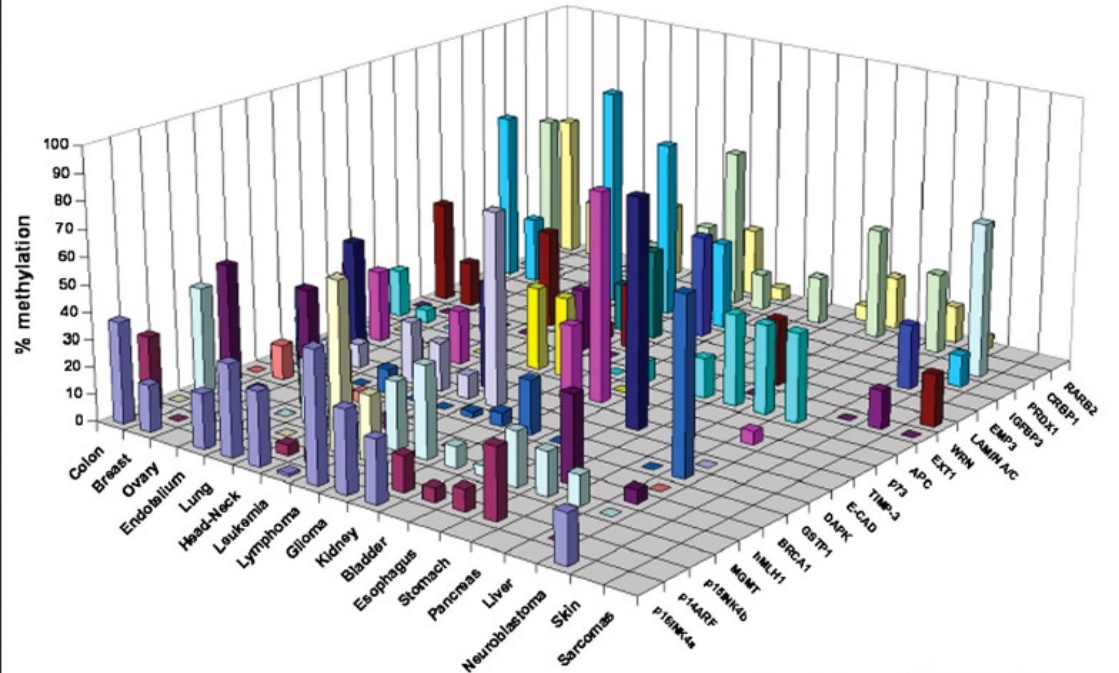


# Graphs Gone Awry

Sometimes, you'll see graphs that are mostly incomprehensible



A CpG Island Hypermethylation Profile of Human Cancer



# Graphs Gone Awry

Or, the relationship between the underlying data and the visual presentation is *really stretched*

## Peak time for sports and leisure

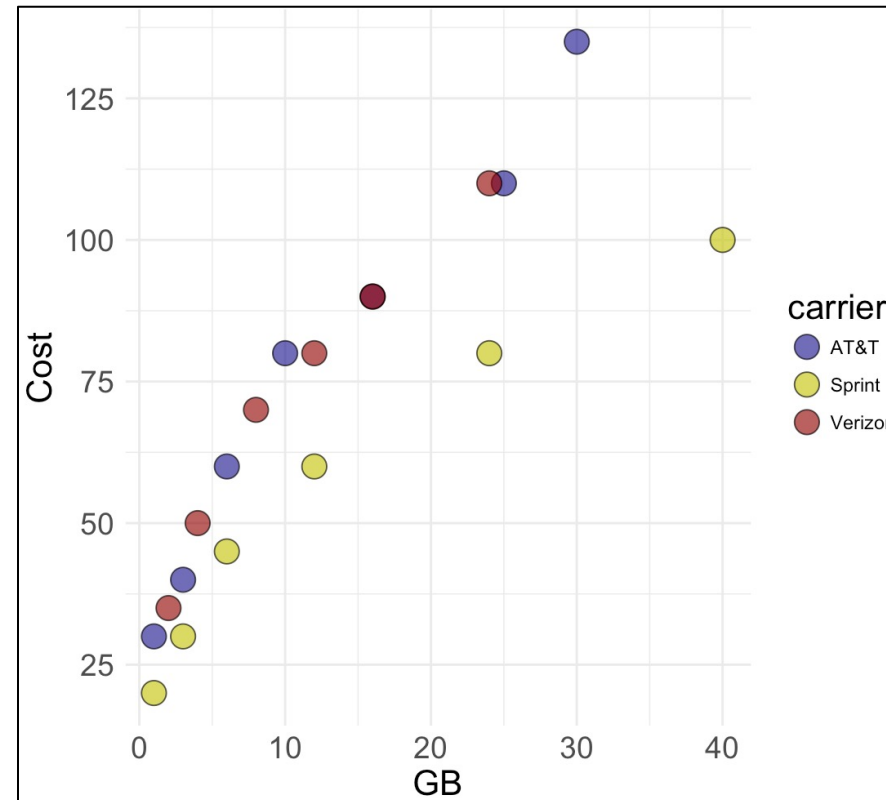
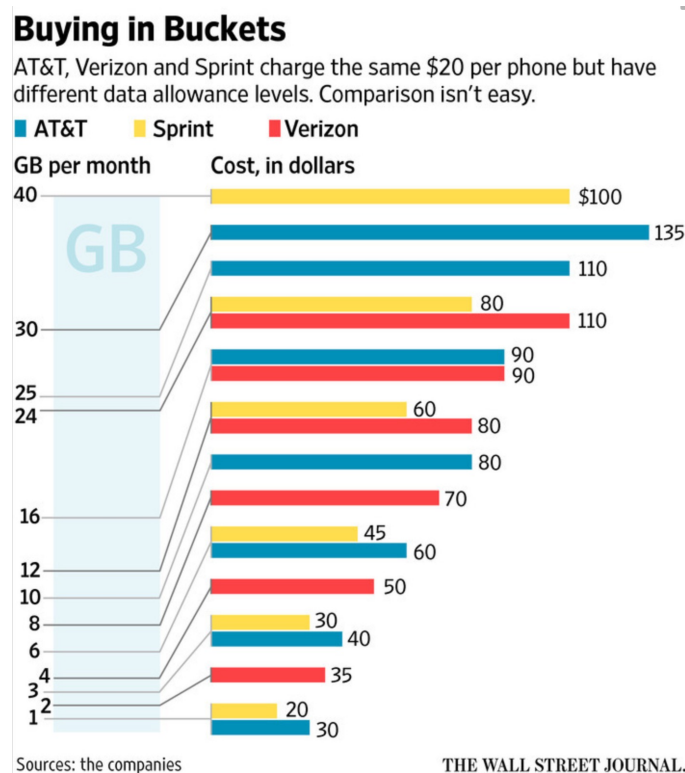
@hnrkIndbrg | Source: American Time Use Survey





# Graphs Gone Awry

Or, the relationship between the underlying data and the visual presentation is *really stretched*



# Graphs Gone Awry

# Graphs Gone Awry

In sum, there are a lot of ways graphs can go haywire.

# Graphs Gone Awry

In sum, there are a lot of ways graphs can go haywire.

Recognizing this in the world around us is important, and avoiding it in our own work even more so!



# What We'll Cover

1. Ways graphs can go wrong (*there are more than you think...*)
2. How to avoid these pitfalls
3. Choosing the right visualization

# What We'll Cover

1. Ways graphs can go wrong (*there are more than you think...*)
2. How to avoid these pitfalls
3. Choosing the right visualization

# Best Practices



# Best Practices

1. Label your axes

# Best Practices

1. Label your axes
2. Follow axis conventions

# Best Practices

1. Label your axes
2. Follow axis conventions
  - Explanatory variable on x axis
  - Simple and intuitive y variable whenever possible

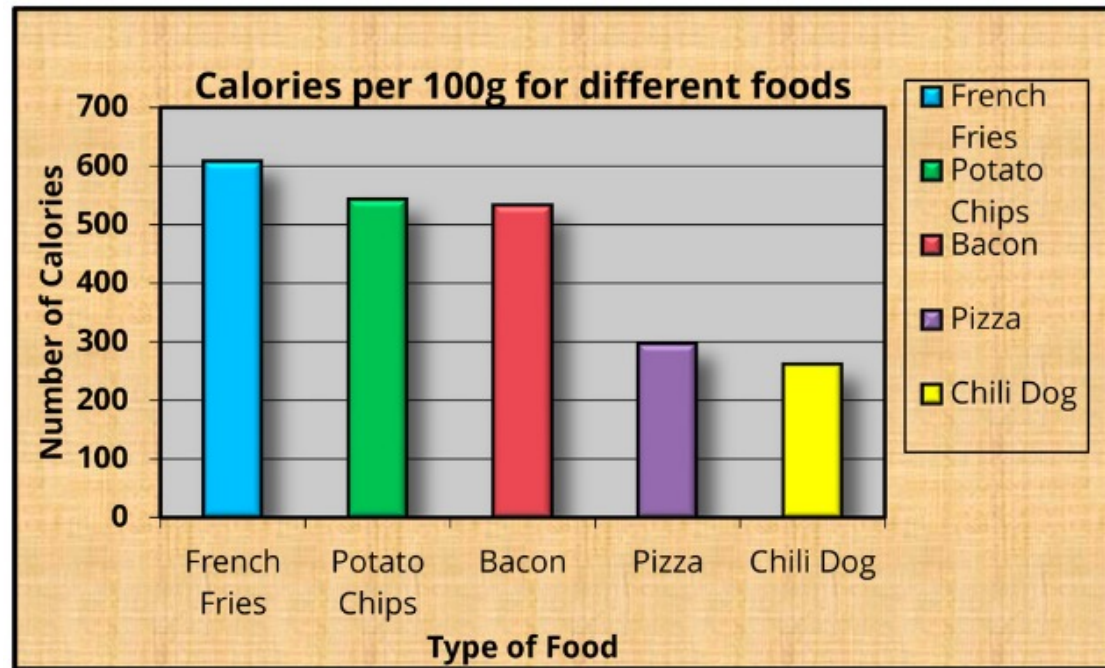
# Best Practices

1. Label your axes
2. Follow axis conventions
  - Explanatory variable on x axis
  - Simple and intuitive y variable whenever possible
3. Include uncertainty (standard error, confidence intervals, or raw data itself)

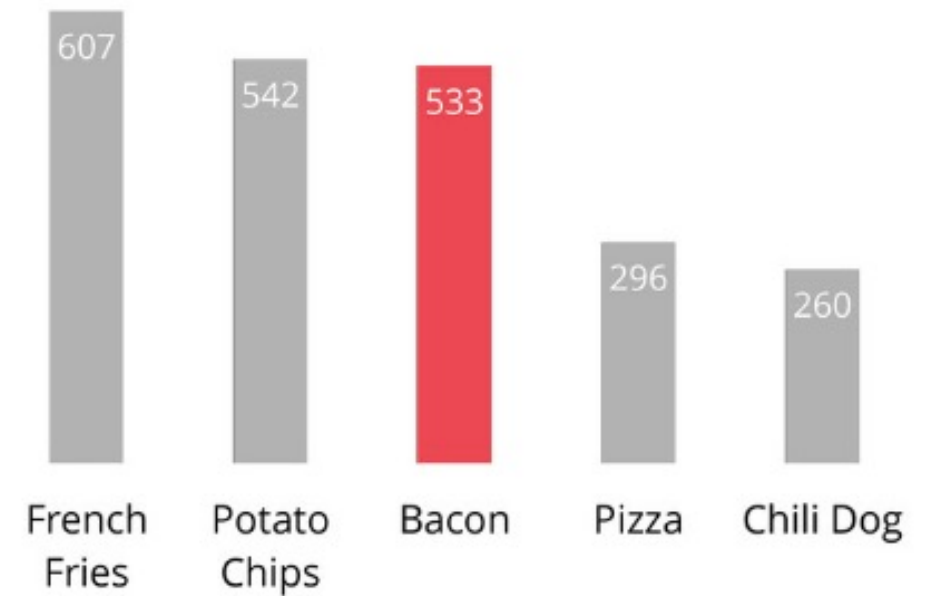
# Best Practices

1. Label your axes
2. Follow axis conventions
  - Explanatory variable on x axis
  - Simple and intuitive y variable whenever possible
3. Include uncertainty (standard error, confidence intervals, or raw data itself)
4. "Remove to improve"

# Best Practices



Calories per 100g



Source: Darkhorse Analytics

# Best Practices

1. Label your axes
2. Follow axis conventions
  - Explanatory variable on x axis
  - Simple and intuitive y variable whenever possible
3. Include uncertainty (standard error, confidence intervals, or raw data itself)
4. "Remove to improve"
5. Choose the right graph for the problem





How to know you're doing it right

# How to know you're doing it right

1. A viewer can easily recognize what's in your graph

# How to know you're doing it right

1. A viewer can easily recognize what's in your graph  
→ clear axis labels and legends, remove to improve

# How to know you're doing it right

1. A viewer can easily recognize what's in your graph  
→ clear axis labels and legends, remove to improve
2. A viewer can easily make the sorts of comparisons you want them to or identify the pattern you're trying to express

# How to know you're doing it right

1. A viewer can easily recognize what's in your graph  
→ clear axis labels and legends, remove to improve
2. A viewer can easily make the sorts of comparisons you want them to or identify the pattern you're trying to express  
→ remove to improve, choose the right graph

# How to know you're doing it right

1. A viewer can easily recognize what's in your graph  
→ clear axis labels and legends, remove to improve
2. A viewer can easily make the sorts of comparisons you want them to or identify the pattern you're trying to express  
→ remove to improve, choose the right graph
3. A viewer can easily figure out what the "point" of the graph is (i.e., what is the key result or finding; don't overdo it!)

# How to know you're doing it right

1. A viewer can easily recognize what's in your graph  
→ clear axis labels and legends, remove to improve
2. A viewer can easily make the sorts of comparisons you want them to or identify the pattern you're trying to express  
→ remove to improve, choose the right graph
3. A viewer can easily figure out what the "point" of the graph is (i.e., what is the key result or finding; don't overdo it!)  
→ Remove to improve, choose the right graph

# How to know you're doing it right

1. A viewer can easily recognize what's in your graph  
→ clear axis labels and legends, remove to improve
2. A viewer can easily make the sorts of comparisons you want them to or identify the pattern you're trying to express  
→ remove to improve, choose the right graph
3. A viewer can easily figure out what the "point" of the graph is (i.e., what is the key result or finding; don't overdo it!)  
→ Remove to improve, choose the right graph
4. A viewer can do all the above with the right level of precision or certainty



# How to know you're doing it right

1. A viewer can easily recognize what's in your graph  
→ clear axis labels and legends, remove to improve
2. A viewer can easily make the sorts of comparisons you want them to or identify the pattern you're trying to express  
→ remove to improve, choose the right graph
3. A viewer can easily figure out what the "point" of the graph is (i.e., what is the key result or finding; don't overdo it!)  
→ Remove to improve, choose the right graph
4. A viewer can do all the above with the right level of precision or certainty  
→ Uncertainty or raw data



# What We'll Cover

1. Ways graphs can go wrong (*there are more than you think...*)
2. How to avoid these pitfalls
3. Choosing the right visualization

# What We'll Cover

1. Ways graphs can go wrong (*there are more than you think...*)
2. How to avoid these pitfalls
3. Choosing the right visualization


# Choosing the right graph for the job

# Choosing the right graph for the job

## Python For Data Science Cheat Sheet

### Seaborn

Learn Data Science Interactively at [www.DataCamp.com](https://www.datacamp.com)



### Statistical Data Visualization With Seaborn

The Python visualization library Seaborn is based on matplotlib and provides a high-level interface for drawing attractive statistical graphics.

Make use of the following aliases to import the libraries:

```
>>> import matplotlib.pyplot as plt
>>> import seaborn as sns
```

The basic steps to creating plots with Seaborn are:

1. Prepare some data
2. Control figure aesthetics
3. Plot with Seaborn
4. Further customize your plot

```
>>> import matplotlib.pyplot as plt
>>> import seaborn as sns
>>> tips = sns.load_dataset("tips")
>>> sns.set_style("whitegrid")
>>> g = sns.lmplot(x="tip", y="total_bill", data=tips, aspect=2)
>>> g = (g.set_axis_labels("Tip", "Total bill (USD)")).set(xlim=(0,10), ylim=(0,100))
>>> plt.title("title")
>>> plt.show(g)
```

**1 Data** Also see Lists, NumPy & Pandas

```
>>> import pandas as pd
>>> import numpy as np
>>> uniform_data = np.random.rand(10, 12)
>>> data = pd.DataFrame({'x': np.arange(1, 101), 'y': np.random.normal(0, 4, 100)})
```

Seaborn also offers built-in data sets:

```
>>> titanic = sns.load_dataset("titanic")
>>> iris = sns.load_dataset("iris")
```

**2 Figure Aesthetics** Also see Matplotlib

**Context Functions**

```
>>> f, ax = plt.subplots(figsize=(5, 6))
>>> sns.set_context("talk")
>>> sns.set_context("notebook", font_scale=1.5, rc={"lines.linewidth": 2.5})
```

**Seaborn styles**

```
>>> sns.set()
>>> sns.set_style("whitegrid")
>>> sns.set_style("ticks", {"tick.major.size": 8, "tick.minor.size": 4})
>>> sns.axes_style("whitegrid")
```

(Re)set the seaborn default  
Set the matplotlib parameters  
Set the matplotlib parameters  
Return a dict of params or use with with to temporarily set the style

**Color Palette**

```
>>> sns.set_palette("husl", 3)
>>> sns.set_palette("magma")
>>> sns.set_palette("magma", "#999999", "#999999", "#999999", "#999999", "#999999")
>>> sns.set_palette(sns.cubehelix_palette(10, 0.7, 0.6))
```

Define the color palette  
Use with with to temporarily set palette  
Set your own color palette

**3 Plotting With Seaborn**

**Axis Grids**

```
>>> g = sns.FacetGrid(titanic, col="survived", row="sex")
>>> g = g.map(plt.hist, "age")
>>> sns.factorplot(x="pclass", y="survived", hue="sex", data=titanic)
>>> sns.lmplot(x="sepal_width", y="sepal_length", hue="species", data=iris)
```

Subplot grid for plotting conditional relationships  
Draw a categorical plot onto a FacetGrid  
Plot data and regression model fits across a FacetGrid

```
>>> h = sns.FacetGrid(iris)
>>> h = h.map(plt.scatter)
>>> sns.pairplot(iris)
>>> i = sns.JointGrid(x="sepal_width", y="sepal_length", data=iris)
>>> i = i.plot(sns.regplot, sns.distplot)
>>> sns.jointplot(x="sepal_width", y="sepal_length", data=iris, kind="kde")
```

Subplot grid for plotting pairwise relationships  
Plot pairwise bivariate distributions  
Grid for bivariate plot with marginal univariate plots  
Plot bivariate distribution

**Categorical Plots**

**Scatterplot**

```
>>> sns.stripplot(x="species", y="petal_length", data=iris)
>>> sns.swarmplot(x="species", y="petal_length", data=iris)
```

Scatterplot with one categorical variable  
Categorical scatterplot with non-overlapping points

**Bar Chart**

```
>>> sns.barplot(x="sex", y="survived", hue="class", data=titanic)
```

Show point estimates and confidence intervals with scatterplot glyphs

**Count Plot**

```
>>> sns.countplot(x="deck", data=titanic, palette="Greens_d")
```

Show count of observations

**Point Plot**

```
>>> sns.pointplot(x="class", y="survived", hue="sex", data=titanic, palette={"male": "g", "female": "m"}, markers=["^", "o"], linestyles=["-", "--"])
```

Show point estimates and confidence intervals as rectangular bars

**Boxplot**

```
>>> sns.boxplot(x="alive", y="age", hue="adult_male", data=titanic)
>>> sns.boxplot(data=iris, orient="h")
```

Boxplot  
Boxplot with wide-form data

**Violinplot**

```
>>> sns.violinplot(x="age", y="sex", hue="survived", data=titanic)
```

Violin plot

**Regression Plots**

```
>>> sns.regplot(x="sepal_width", y="sepal_length", data=iris, ax=ax)
```

Plot data and a linear regression model fit

**Distribution Plots**

```
>>> plot = sns.distplot(data.y, kde=False, color="b")
```

Plot univariate distribution

**Matrix Plots**

```
>>> sns.heatmap(uniform_data, vmin=0, vmax=1)
```

Heatmap

**4 Further Customizations** Also see Matplotlib

**Axisgrid Objects**

```
>>> g.despine(left=True)
>>> g.set_ylabels("Survived")
>>> g.set_xticklabels(rotation=45)
>>> g.set_axis_labels("Survived", "Sex")
>>> h.set(xlim=(0, 5), ylim=(0, 5), xticks=[0, 2.5, 5], yticks=[0, 2.5, 5])
```

Remove left spine  
Set the labels of the y-axis  
Set the tick labels for x  
Set the axis labels  
Set the limit and ticks of the x-and y-axis

**Plot**

```
>>> plt.title("A Title")
>>> plt.ylabel("Survived")
>>> plt.xlabel("Sex")
>>> plt.ylim(0, 100)
>>> plt.xlim(0, 10)
>>> plt.setp(ax, yticks=[0, 5])
>>> plt.tight_layout()
```

Add plot title  
Adjust the label of the y-axis  
Adjust the label of the x-axis  
Adjust the limits of the y-axis  
Adjust the limits of the x-axis  
Adjust a plot property  
Adjust subplot params

**5 Show or Save Plot** Also see Matplotlib

```
>>> plt.show()
>>> plt.savefig("foo.png")
>>> plt.savefig("foo.png", transparent=True)
```

Show the plot  
Save the plot as a figure  
Save transparent figure

**Close & Clear** Also see Matplotlib

```
>>> plt.cla()
>>> plt.clf()
>>> plt.close()
```

Clear an axis  
Clear an entire figure  
Close a window

**DataCamp**  
Learn Python for Data Science Interactively

Source: <https://www.kaggle.com/getting-started/126958>

# Choosing the right graph for the job

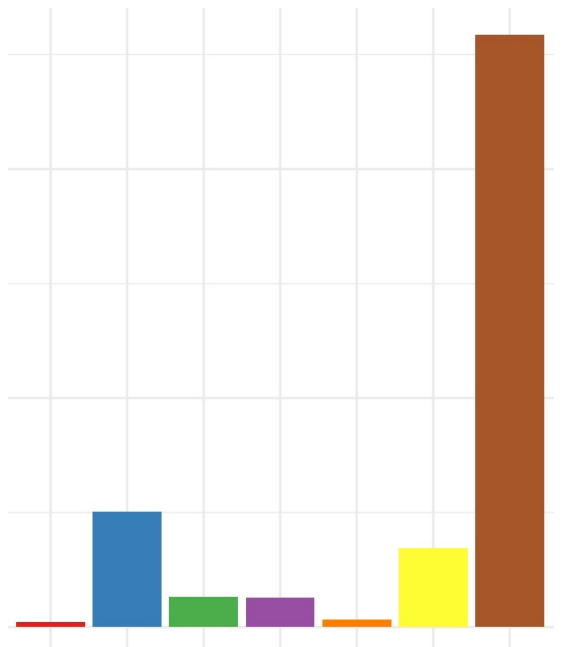
What kind of variables are we working with?

- Categorical / discrete
- Numerical / continuous

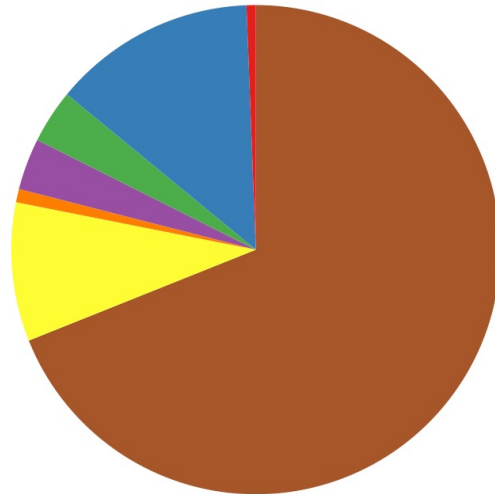
Which variable(s) are explanatory and which are being explained (response variable)?

# Choosing the right graph for the job

Categorical response variable, no explanatory variable



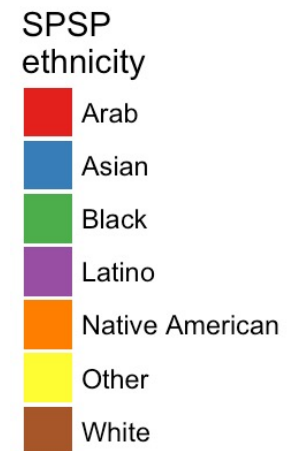
Histogram barplot  
(counts or proportions)



Pie chart



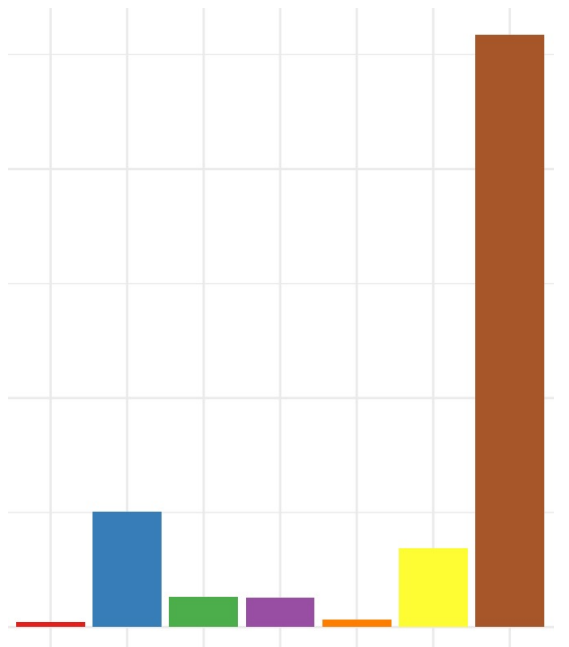
Stacked barplot



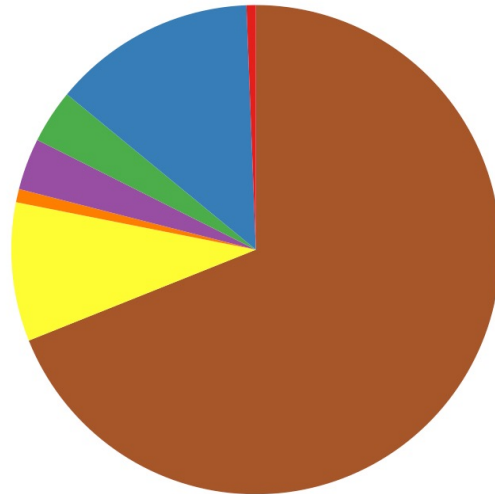


# Choosing the right graph for the job

Categorical response variable, no explanatory variable



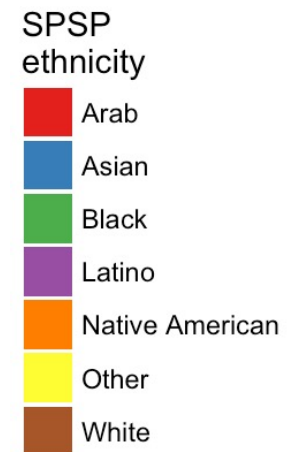
Histogram barplot  
(counts or proportions)



Pie chart

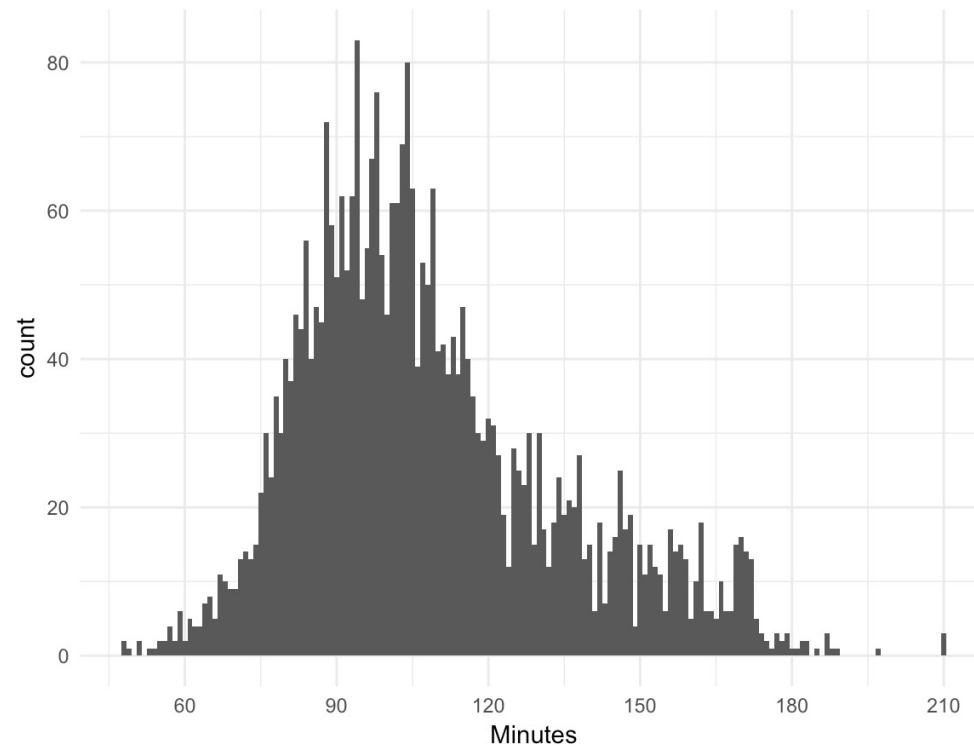


Stacked barplot

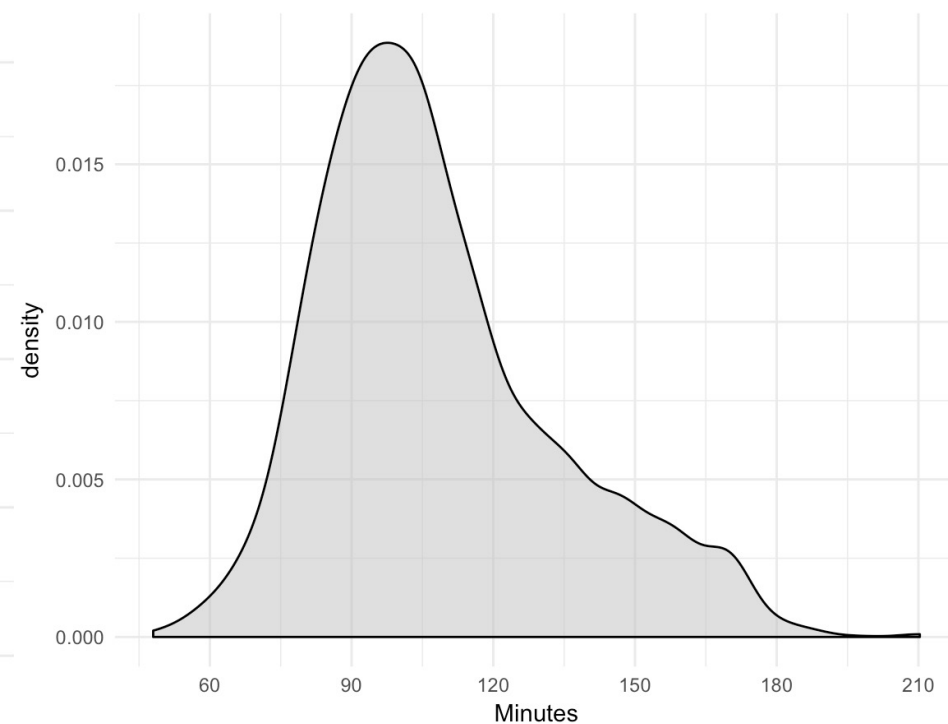


# Choosing the right graph for the job

Numerical response variable, no explanatory variable



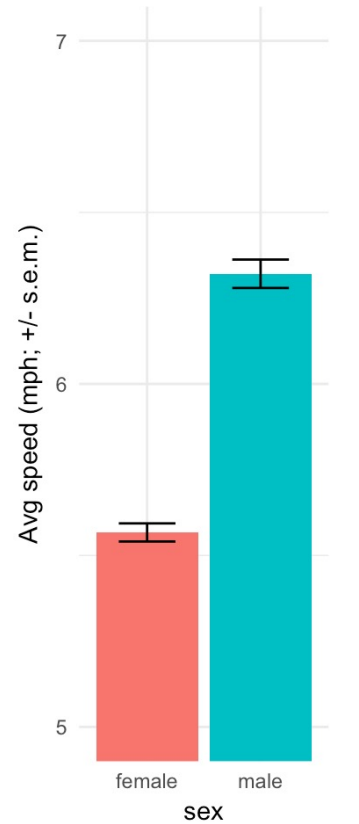
Histogram!



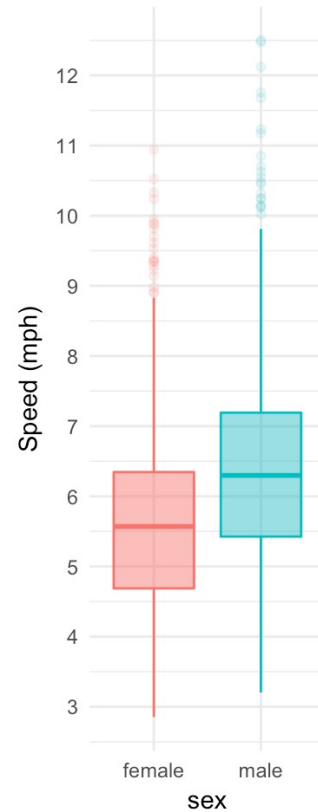
Smoothed density

# Choosing the right graph for the job

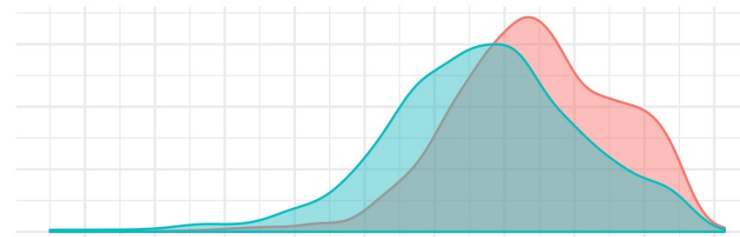
Numerical response variable, categorical explanatory variable



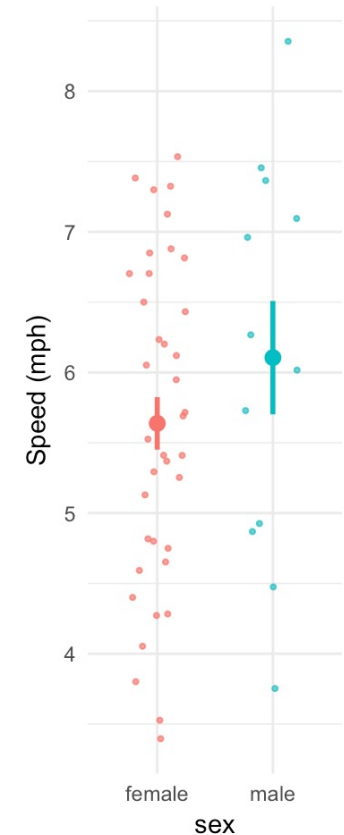
Mean + error



Boxplot

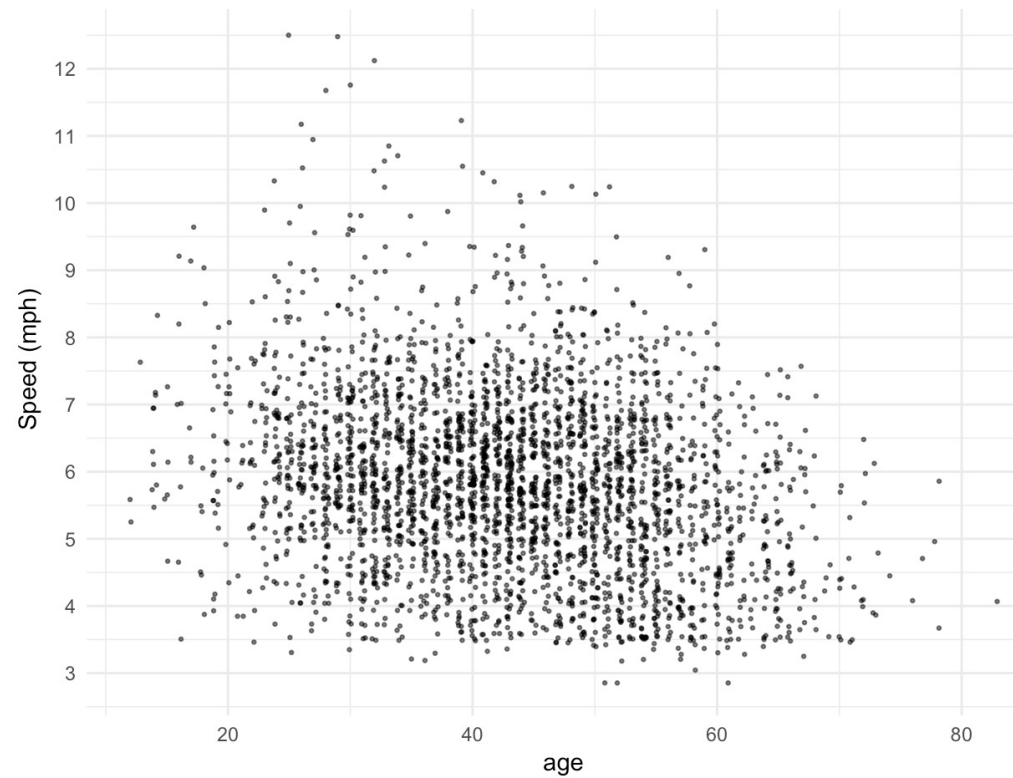


Densities (or histograms)



# Choosing the right graph for the job

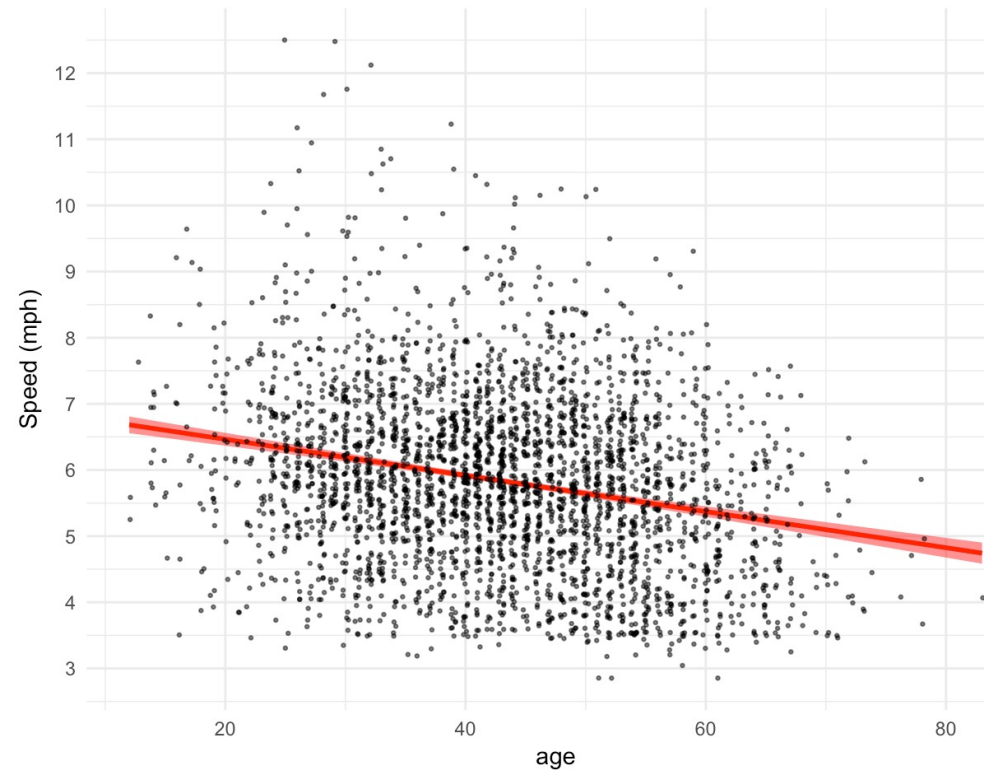
Numerical response variable, numerical explanatory variable



Scatterplot

# Choosing the right graph for the job

Numerical response variable, numerical explanatory variable

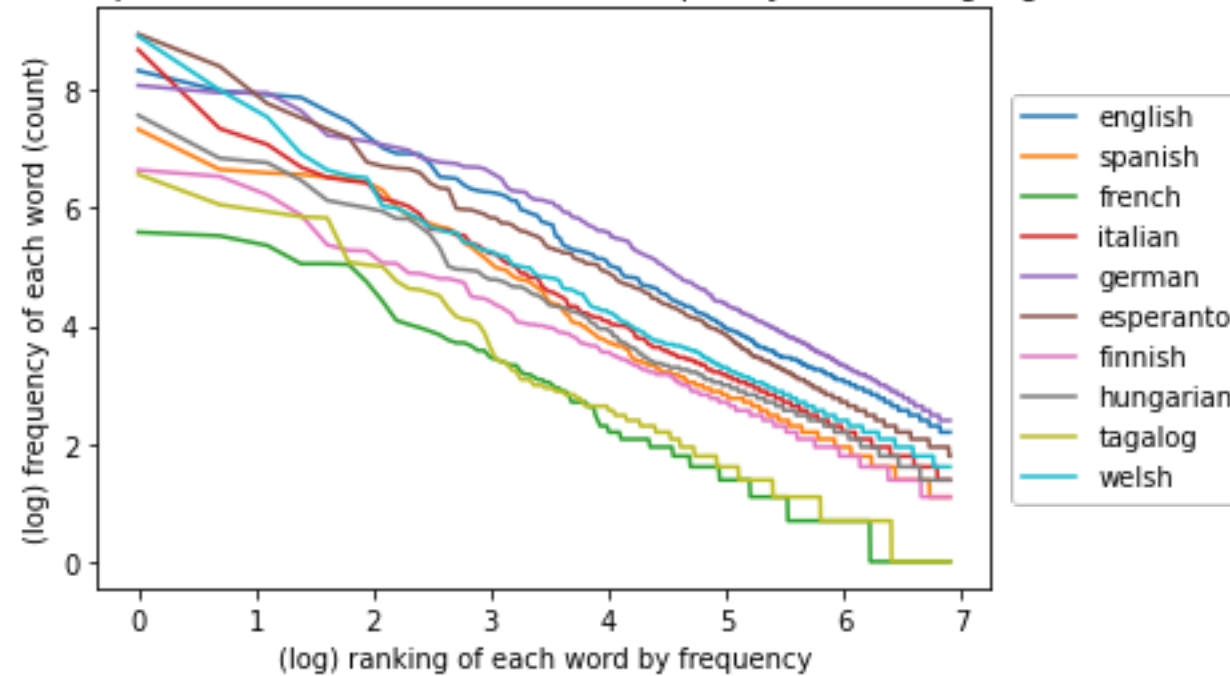


Include line of best fit!

# Choosing the right graph for the job

Numerical response variable, numerical *and* categorical explanatory variables

Relationship between word rank and word frequency across language texts



More options here, but this is usually a good start!

# Choosing the right graph for the job

- This obviously isn't **everything** when it comes to visualization options
- But this should get you on the right path most of the time!





# What We'll Cover

1. Ways graphs can go wrong (*there are more than you think...*)
2. How to avoid these pitfalls
3. Choosing the right visualization

