# Intro to Data Science in p5.js

Adam Prado

# Goals of Unit (CS, Stats, Civics)

- Further developing p5 coding skills.

- Load, access, and manipulate date with p5.

- Practice with visualizing data.

- Making conclusions based on examining data.

- Checking for errors and bias in data or conclusions.

- Create something to improve yourself or your community.

- Brief exposure to "Data Science"

# Unit Overview

1. Intro to Data Science + importing data into p5.

2. Creating functions with univariate data

* Cleaning and filtering data sets

3. Bar Graph (map () )

4. Circle Graph

5. Bivariate Data + Line Graph

6. Scatter Plot

7. Linear regression (causation v. correlation

8. Analyzing results and making conclusions

*. Big Data, Machine Learning,

9. Concerns about bias
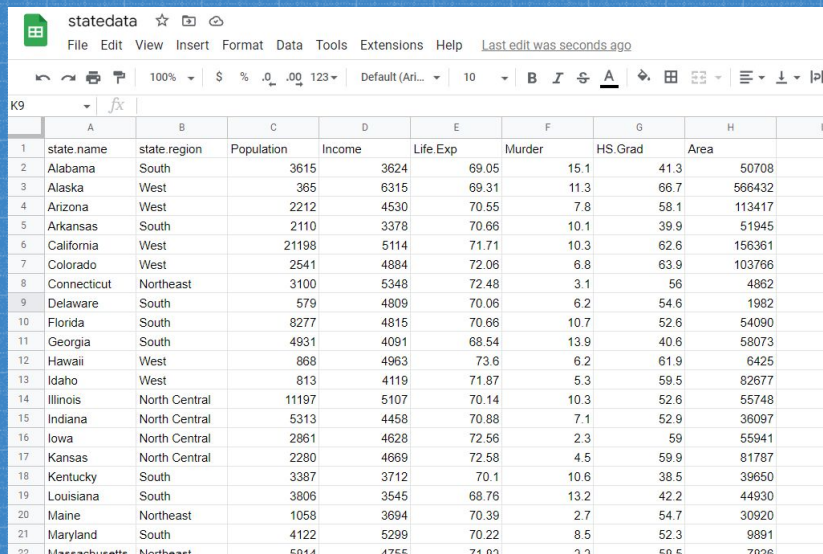
10-12. Data Science group projects

# Standards: Data Analysis and Visualization

9-12.CT.2 Collect and evaluate data from multiple sources for use in a computational artifact.

9-12.CT.3 Refine and visualize complex data sets to tell different stories with the same data set.

9-12.CT.7 Design or remix a program that utilizes a data structure to maintain changes to related pieces of data.

# 1) Intro to Data Science + importing data into p5.



- What is data science?

- File types (pdf, txt, json, xls, csv)

- Spreadsheet to .csv

- loading csv into p5.js

- preload()

- loadTable()

- data.columns

- data.getRow()

▶  ■   ☐ Auto-refresh    Vaulted spider ✎

| Sketch Files ▾ | ‹ sketch.js● | Preview |

Sketch Files ▾

📄 index.html
📄 sketch.js
📄 statedata.csv
📄 style.css

```javascript
1   let data;
2
3   function preload(){
4       data = loadTable("statedata.csv","csv","header")
5
6   }
7
8   function setup() {
9       createCanvas(400, 400);
10      background(220);
11
12      let categories = data.columns
13      console.log(categories)
14
15      let num = data.getRowCount()
16      console.log(num)
17
18      let stNames = data.getColumn(0)
19      console.log(stNames)
20  }
21
22  function draw() {
23      //background(220);
```

Console                                          Clear ∨

"HSGrad", "Area"]


50

▶ (50) ["Alabama", "Alaska", "Arizona", "Arkansas", "California", "Colorado", "Connecticut", "Delaware", "Florida", "Georgia", …]

# 2) Creating functions with univariate data

- Discussion of Univariate Data

- Starter code with (Length of fish) and skeleton of comments with functions, instructions, and test conditions.

- Function to find mean

- Creating code to count values that meet a given condition

    - # fish over a certain size

    - # of fish between a range of 2 values

*) Function to find median, including .sort()

▶   ■   ☐ Auto-refresh   Mean Median Count ✎   by ajprado@gmail.com

> sketch.js •                                                    Saved: 3 days ago

```
 9    createCanvas(400, 400);
10    background(220);
11
12    let lengths = data.getColumn(0)
13    //print(lengths)
14    print("mean fish length: " + mean(lengths))
15    print("median fish length: " + median(lengths))
16    let countBig = 0
17▾   for(let i=0;i<lengths.length;i++){
18      if(lengths[i]>30){
19        countBig++
20      }
21    }
22    print("number of Big Fish is: "+ countBig)
23
24    let countBetween = 0
25▾   for(let i=0;i<lengths.length;i++){
26      if(lengths[i]>20 && lengths[i]<30){
27        countBetween++
28      }
29    }
30    print("number of fish between 20 and 30: "+ countBetween)
31
32  }
```

Console                                                              Clear

```
mean fish length: 26.02857142857143
median fish length: 27
number of Big Fish is: 9
number of fish between 20 and 30: 12
```

>

# *. Cleaning and filtering data sets

- What are some of the problems with collecting data?

- What does cleaning data mean?

- Using a spreadsheet to look over data

- Changing heading names

- Removing entries with missing/erroneous values.

- Formatting data to work best in p5.js

# 3. Bar Graph (map () )



```
map()

Examples

let value = 25;
let m = map(value, 0, 100, 0, width);
ellipse(m, 50, 10, 10);
```

- Starter code with examples of using p5 map() function

- Practice using map() to edit the scale of numbers

- Starter code with data of animals and "cuteness" rating. (possibly scaffolded with axis set up)

- Use map to convert cuteness from 0 to 10 to a scale that will look better on the canvas.

*) Color the bars of the graph differently depending on the number of legs the animal has.

```
23        fill("blue")
24      }else if (legs[i]==2){
25        fill("green")
26      }else{
27        fill("red")
28      }
29
30
31      let mapCute = map(cute[i],0,10,0,200)
32      rect(110,90+(i*20),mapCute,10)
33    }
34    fill(0)
35    line(110,70,110,230)
36    line(110,230, 350, 230)
37    for(let i=0; i<11;i++){
38      textAlign(CENTER);
39      text(i, map(i,0,10,110,310),250)
40    }
41    text("Cuteness",220,270)
42  }
43
44  function draw() {
45    //background(220);
46
47  }
```

Console                                                Clear ⌄

▶ (7) ["9", "8", "1", "2", "10", "7", " 3"]

# 4. Circle Graph

- Starter code with racial demographics of US.

- map the % for each race from 0 to 2PI to get the angle for each sector(brief explanation of radian)

- use arc() with angle measures to construct circle graph. (more technical so likely code-along)
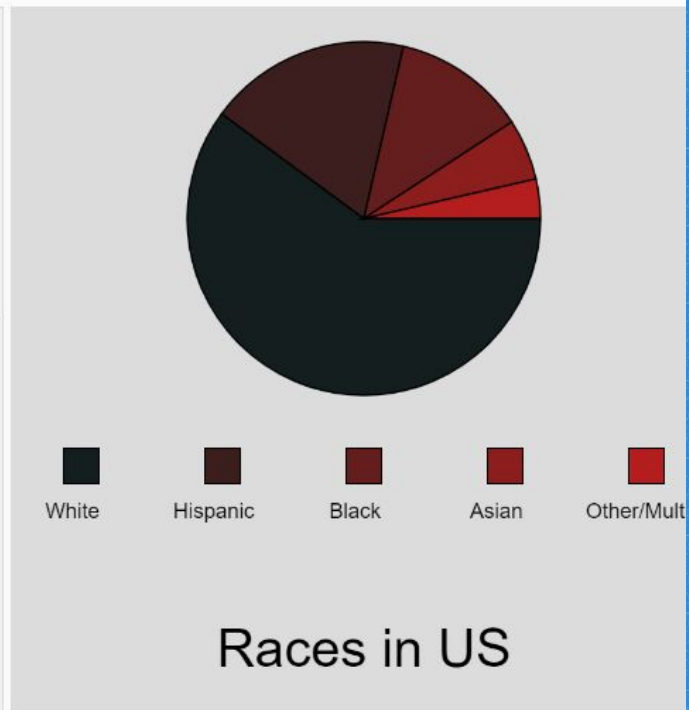
- Create title and set up key to label each section

```javascript
16    let angles = []
17    for (let i=0; i<percent.length;i++){
18      angles.push(map(percent[i],0,100,0,2*PI))
19    }
20    let color= 20;
21    let start = 0
22    let stop;
23    for(let i=0; i<angles.length; i++){
24      stop = start + angles[i]
25      fill(color,30,30)
26      arc(200,120,200,200,start,stop,PIE)
27      color += 40
28      start = stop;
29      rect(30+(i*80),250,20,20)
30      fill(0)
31      textAlign(CENTER,TOP)
32      textSize(12)
33      text(race[i],35+(i*80),280)
34    }
35    textSize(30)
36    text("Races in US", 200,350)
37    print(angles)
38  }
39
40  function draw() {
```

Races in US

White    Hispanic    Black    Asian    Other/Mult

# 5. Bivariate Data + Line Graph

- What is bivariate data?

- How can we use data pairs to show change over time?

- Starter code with data for temperature difference from the mean for each month since 1880

- Code comments leading students to create the line graph.

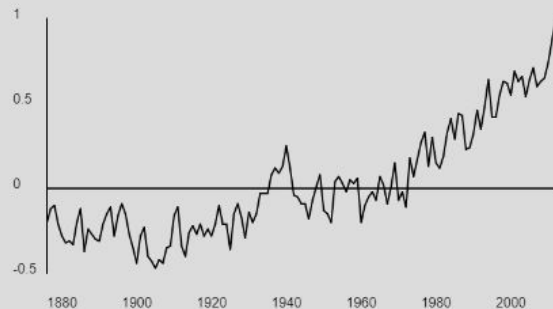- Use line() with map() to connect consecutive points (scaffolded hints for how to set up each map() scale conversion.

```
10    background(220);
11    let numRows = data.getRowCount();
12    let year = data.getColumn('Year')
13    let tempDiff = data.getColumn('Difference from mean')
14    let yearScaled = []
15    let tempScaled = []
16    for(let i=0; i<year.length; i++){
17      yearScaled.push(map(year[i],1880,2016,50,350))
18      tempScaled.push(map(tempDiff[i],-0.5,1,250,100))
19    }
20    line(50,200,350,200)
21    line(50,100,50,250)
22    for(let i=0; i<year.length-1; i++){
23
    line(yearScaled[i],tempScaled[i],yearScaled[i+1],tempScaled[i+1])
24    }
25    textSize(8)
26    for(let i = 1880; i<2020;i +=20){
27      text(i, map(i,1880,2016,50,350),270)
28    }
29    for(let i = -0.5; i<1.5;i +=0.5){
30      text(i, 30,map(i,-0.5,1,250,100))
31    }
32    textSize(12)
33    text("Average global surface temperature difference from the
```

Average global surface temperature difference from the mean

Console      Clear

## 6. Scatter Plot

- Starter code comparing home prices and age of home.

- Code skeleton with comments leading students through the activity. Similar to line graph but using point() instead of lines.

- Discussion on what information can be found from the scatterplot.

```
12   let age = data.getColumn(2)
13   let price = data.getColumn(4)
14
15   console.log(age)
16   textSize(28)
17   text("Age and Price of Homes",50,50)
18
19   strokeWeight(1)
20   line(98,98,98,350)
21   line(98,350,400,350)
22   textSize(14)
23   text("Price of Home/ sqft",200,400)
24   text("Age of home",10,200)
25   strokeWeight(5)
26   for(let i=0; i<age.length;i++){
27       point(map(price[i],0,1700,100,400),map(age[i],0,45,350,100))
28   }
29   textSize(10)
30   for(let x=0; x<1800;x+=250){
31     text(x, map(x,0,1700,100,400), 370)
32   }
33 }
```

Console     Clear



Age and Price of Homes

Age of home

Price of Home/ sqft

0   250   500   750   1000   1250   1500   1750

## 7. Linear regression

- Starting with a fully running program that allows students to plot points and it creates the trend line.

- Discussion of what is the meaning of the line.

- Strong vs. weak correlation examples

- Negative vs. positive correlation examples

- Given a mostly working data set with a prebuilt regression line calculator and having students finish it up with code comments.

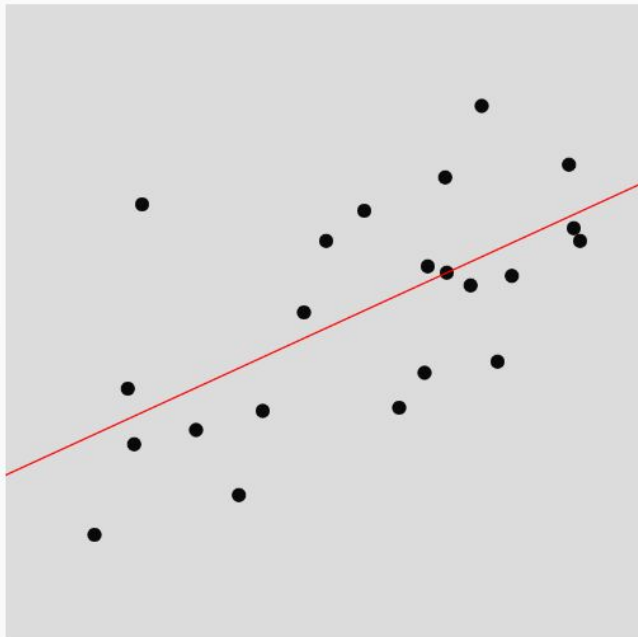> sketch.js                                    Saved: just now        Preview

```javascript
16    background(220);
17    for(var i =0; i<data.length; i++){
18      var x = map(data[i].x, 0, 1, 0, width);
19      var y = map(data[i].y, 0,1,height, 0);
20      fill(10);
21      stroke(10);
22      ellipse(x,y,8,8);
23    }
24
25    if(data.length > 1) {
26      linearRegression();
27      drawLine();
28    }
29  }
30
31  var m = 1;
32  var b = 0;
33
34  function drawLine(){
35    var x1 = 0;
36    var y1 = m * x1 + b;
37    var x2 = 1;
```

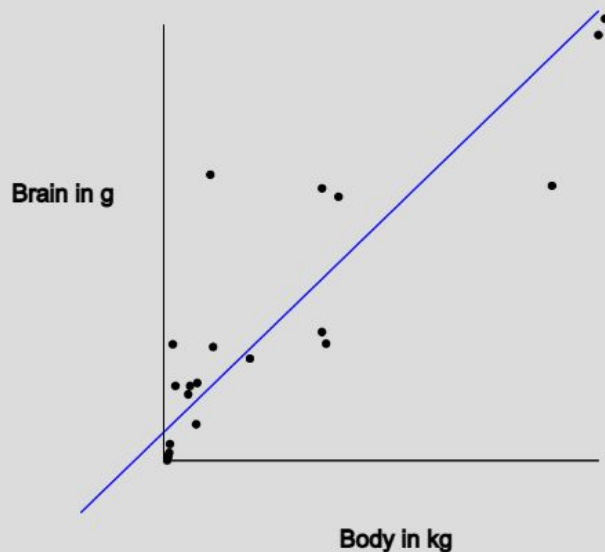Console                                                    Clear  ⌄

```javascript
 8  function setup() {
 9    createCanvas(400, 450);
10    background(220);
11    let numRows = data.getRowCount();
12    let name = data.getColumn(0)
13    let body = data.getColumn(1)
14    let brain = data.getColumn(2)
15    console.log(brain)
16    textSize(28)
17    text("Body/Brain size of animals",50,50)
18
19   stroke("■blue")
20    findLinearReg()
21    stroke("■black")
22    strokeWeight(1)
23    line(98,98,98,350)
24    line(98,350,350,350)
25    textSize(14)
26    text("Body in kg",200,400)
27    text("Brain in g",10,200)
28    strokeWeight(5)
29    for(let i=0; i<name.length;i++){
```


Body/Brain size of animals

Console      Clear
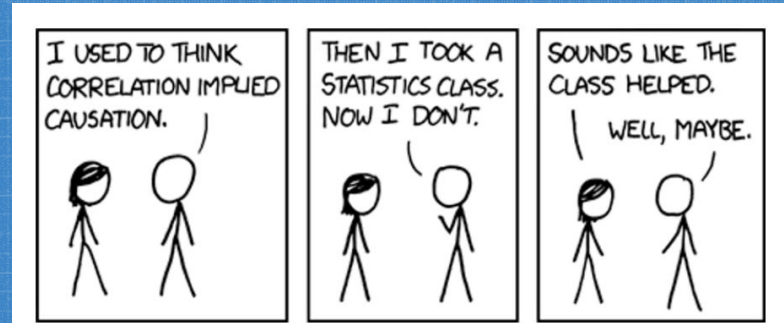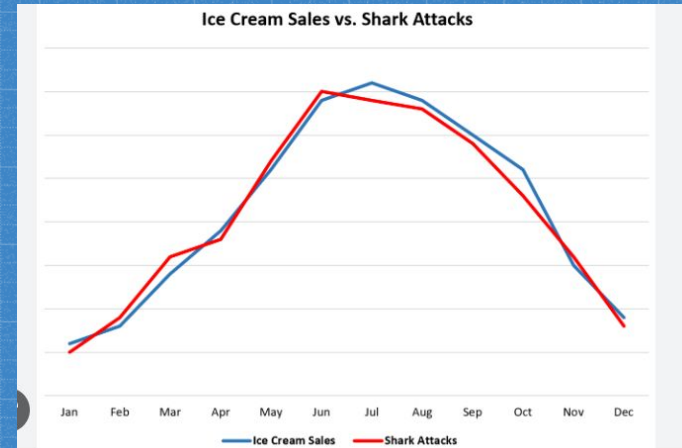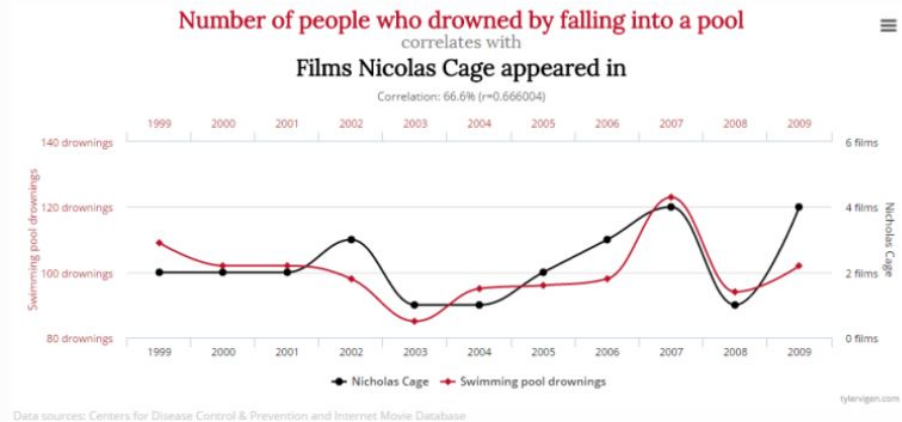
# 8. Analyzing results and making conclusions

- Activity with examples of ridiculous correlated data.  Discussion of the important difference between correlation and causation.

- Example of data conclusions using overgeneralization, discussion the scope of conclusions that can be made.

Number of people who drowned by falling into a pool correlates with Films Nicolas Cage appeared in — Correlation: 66.6% (r=0.666004)


Ice Cream Sales vs. Shark Attacks


TEDxDelft — In Nature there was a study in 1999 that showed

1. Thing A caused Thing B (causality)
2. Thing B caused Thing A (reversed causality)
3. Thing A causes Thing B which then makes Thing A worse (bidirectional causality)
4. Thing A causes Thing X causes Thing Y which ends up causing Thing B (indirect causality)
5. Some other Thing C is causing both A and B (common cause)
6. It's due to chance (spurious or coincidental)

# 9. Concerns about bias

- Video regarding concerns of bias data/conclusions.

- Article with examples of several of the types of statistical bias.

- Discussion of reliability of data (random sample, collection techniques, inclusive data …)

**Types of bias in statistics:**

✔ Confirmation bias
✔ Selection bias
✔ Outlier bias
✔ Observer bias
✔ Funding bias
✔ Omitted variable bias
✔ Survivorship bias

OCTOBER 24, 2020

BLOG, SCIENCE POLICY, SPECIAL EDITION: SCIENCE POLICY AND SOCIAL JUSTICE

**Racial Discrimination in Face Recognition Technology**

# Amazon scraps secret AI recruiting tool that showed bias against women

By Jeffrey Dastin

8 MIN READ

SAN FRANCISCO (Reuters) - Amazon.com Inc's AMZN.O machine-learning specialists uncovered a big problem: their new recruiting engine did not like wo

# 10-12. Data Science group projects

- Find a data set you are interested in exploring, preferably something that impacts you or your community.

- Create at least two data visualizations from it using p5.

- Explain any conclusions you can make about the data including justification.

- Discuss any possible concerns about bias or errors?

- 5 minutes presentation to the class about what you discovered.