

Does Whiskey Cure Diabetes?





Hypothetical Study

A study once found that patrons at a liquor store who exclusively purchased pricier whiskey (> \$40 a bottle) were less likely to have Type II diabetes.

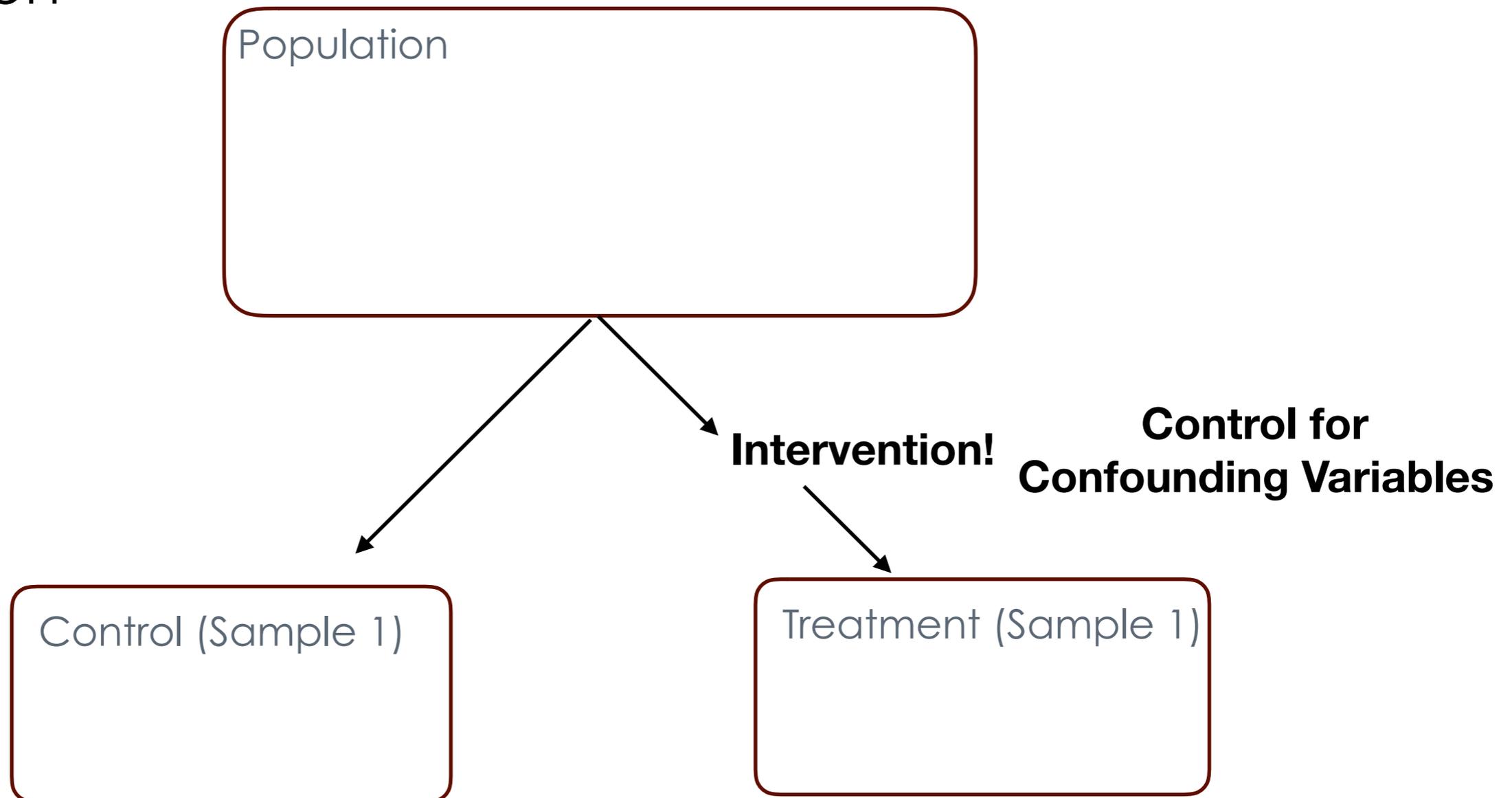
What can we conclude?

- (a) Expensive whiskey mitigates the effects of diabetes
- (b) Cheap whiskey exacerbates the effects of diabetes
- (c) Expensive whiskey drinkers are probably wealthier and have better overall diet and exercise.



Recap: Correlation v.s. Causation

(Last Time) Impossible to guess a causal link without more information



Suggests Causation...



CHIDATA

What if we don't care about causality?

Patrons who exclusively purchased pricier whiskey (> \$40 a bottle).

What can we recommend to a customer?

- (a) Other expensive/luxury products (because they can afford it)
- (b) Whiskey glasses (so they can drink the whiskey)
- (c) Chocolate (because they are less likely to have diabetes)



Degrees of Predictive Value

Patrons who exclusively purchased pricier whiskey (> \$40 a bottle).

Likelihood of following through with a purchase

**More
accurate**



- (a) Other expensive/luxury products (because they can afford it)
- (b) Whiskey glasses (so they can drink the whiskey)
- (c) Chocolate (because they are less likely to have diabetes)



Degrees of Predictive Value

Two variables: X and Y How well does X predict Y ?

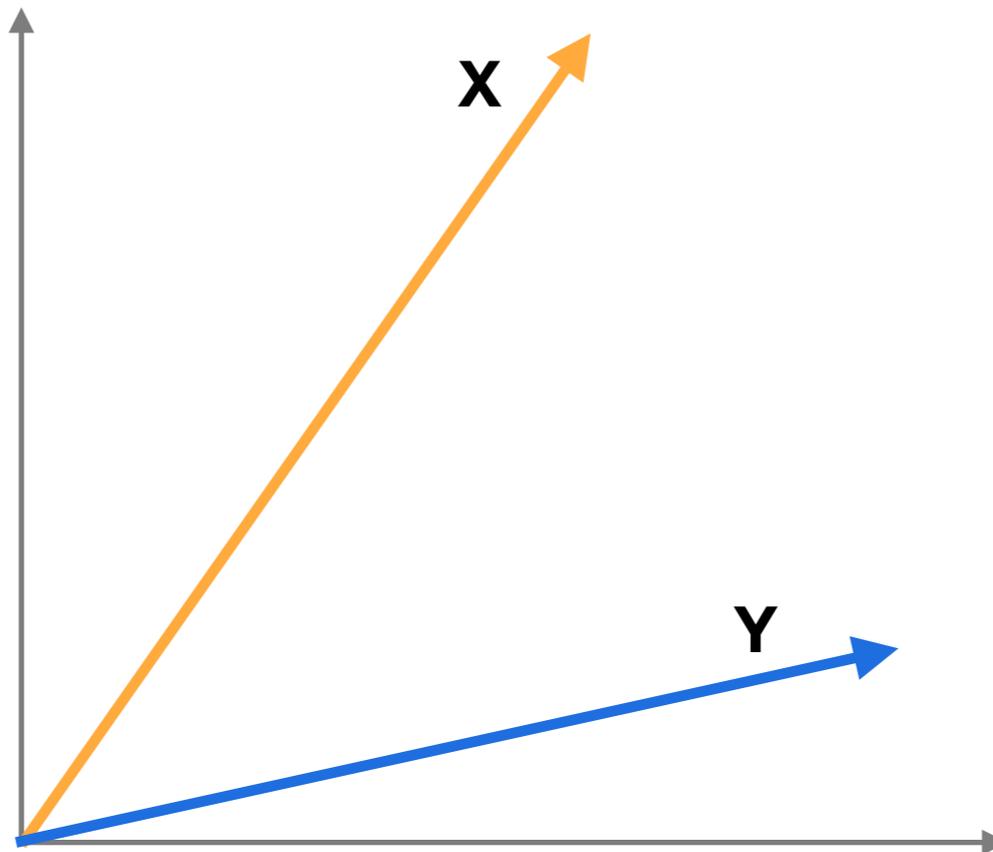
What does this mean and how do we measure predictive value with uncertain quantities?



High School Geometry

Two **vectors**: X and Y

How well does X **align with** Y?



Are they going in the same direction? How closely?

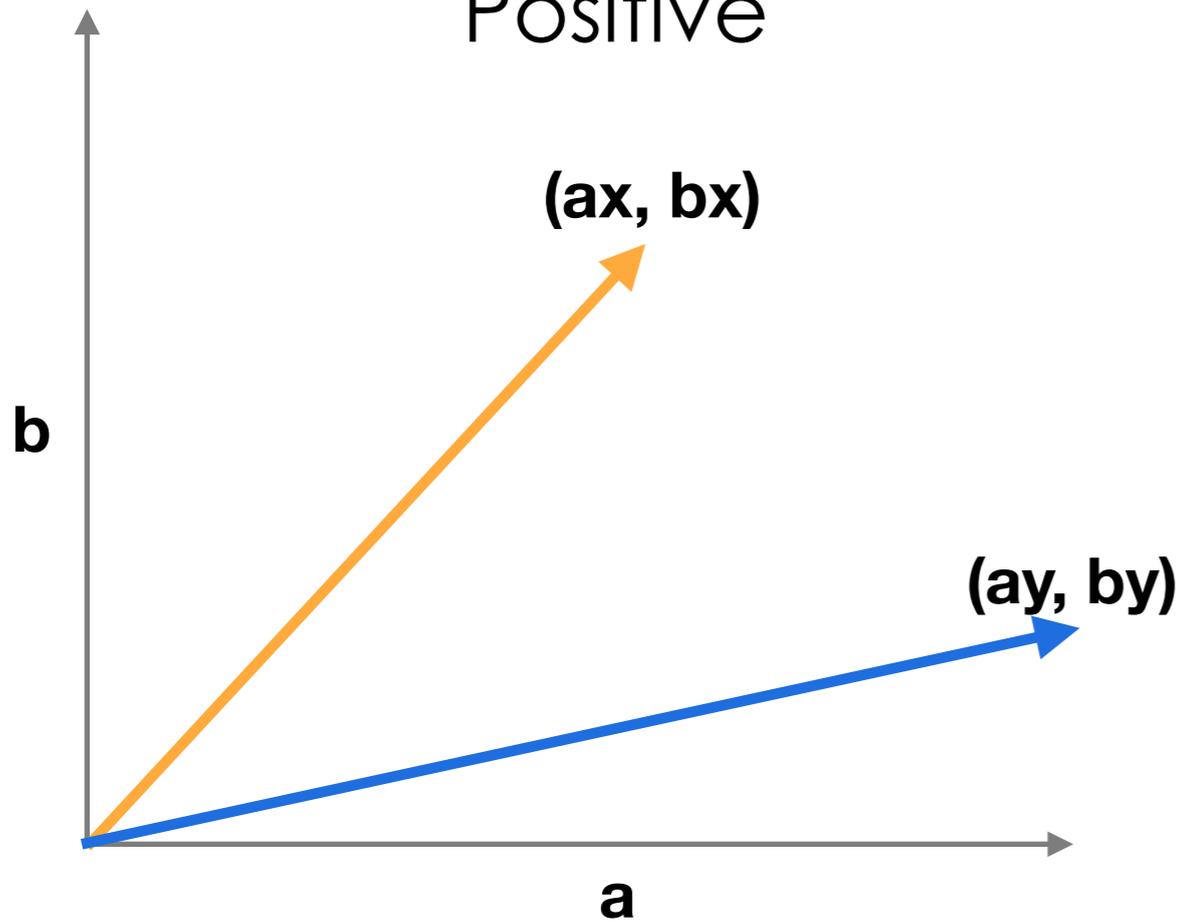


Dot Products

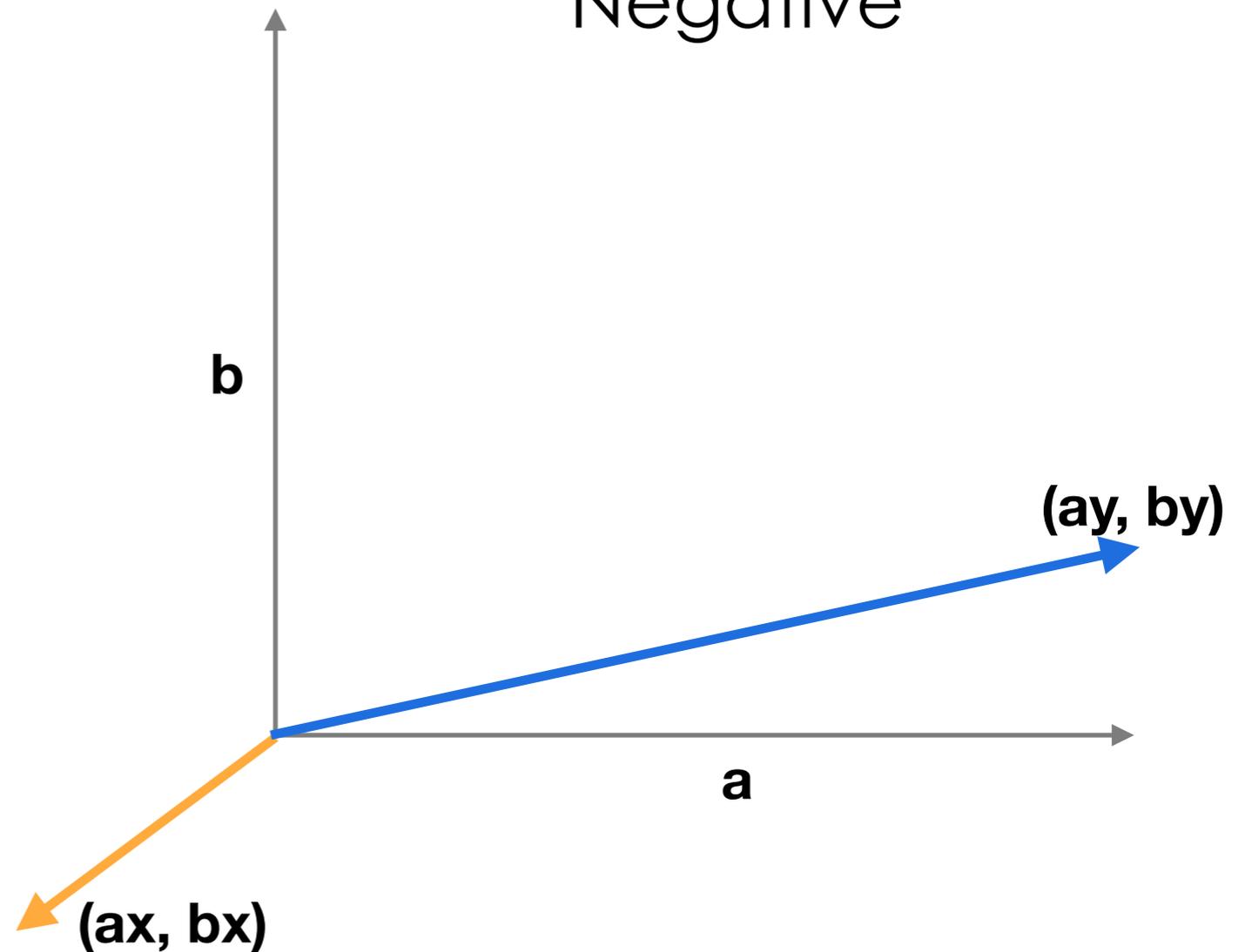
Two **vectors**: X and Y

$$x \cdot y = ax*bx + ay*by$$

Positive



Negative

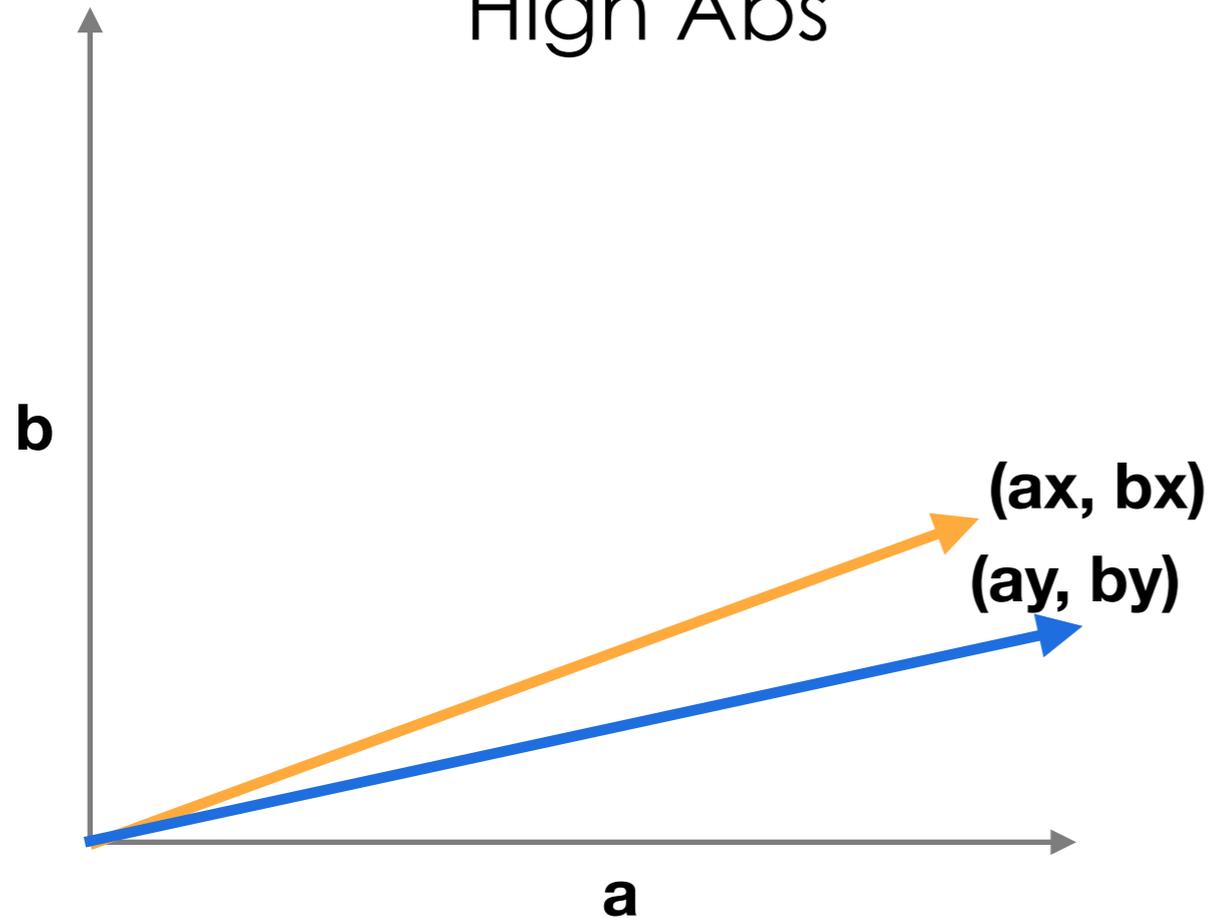


Absolute Value Depends on Angle

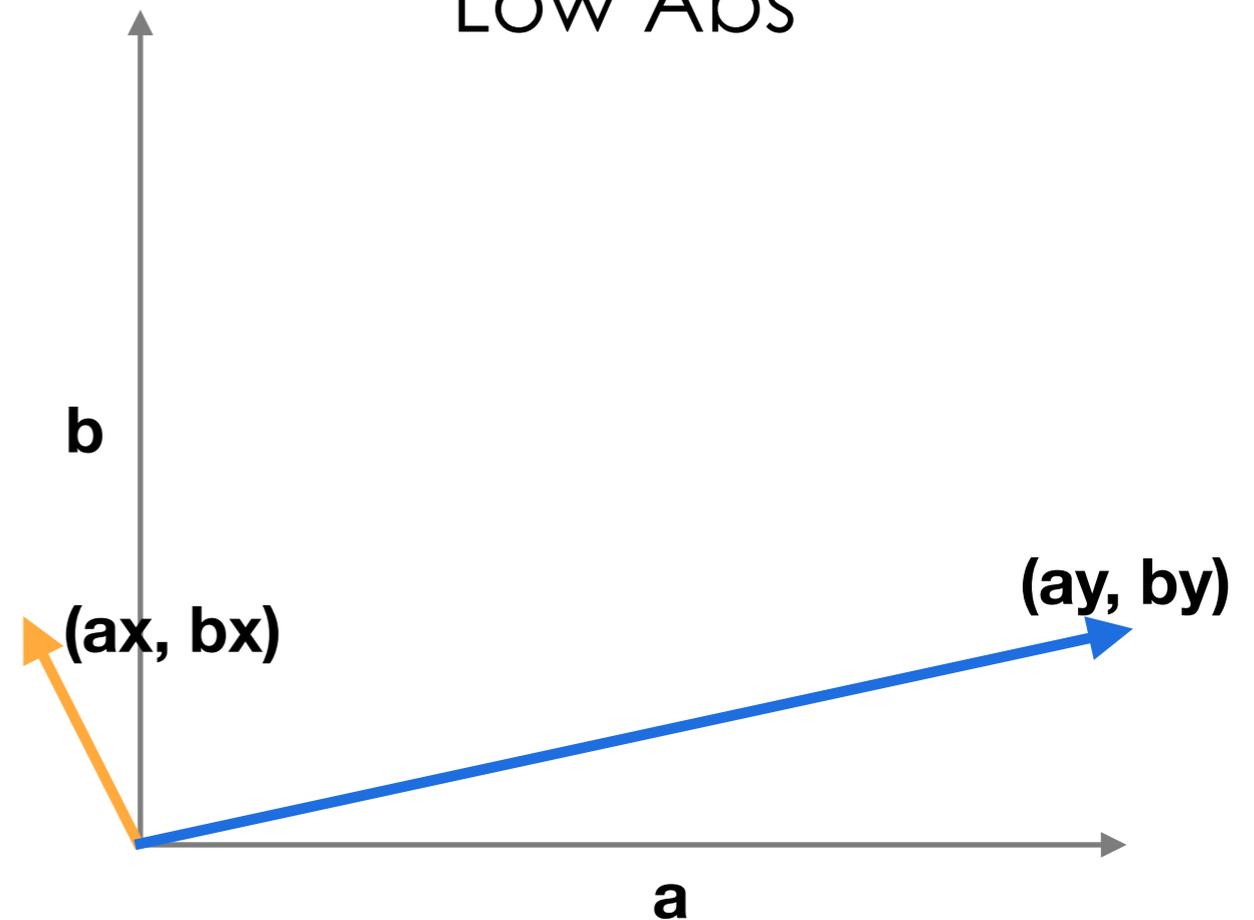
Two **vectors**: X and Y

$$x \cdot y = ax*bx + ay*by$$

High Abs



Low Abs

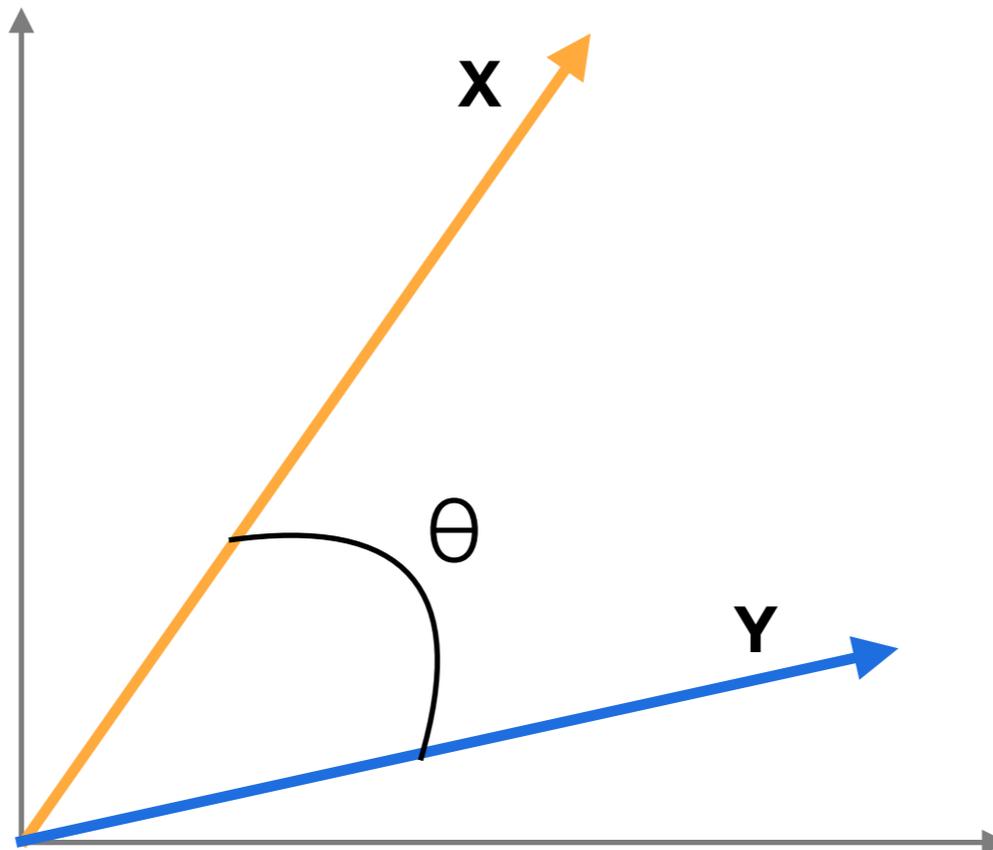




High School Geometry

Two **vectors**: X and Y

Calculate the Angle?



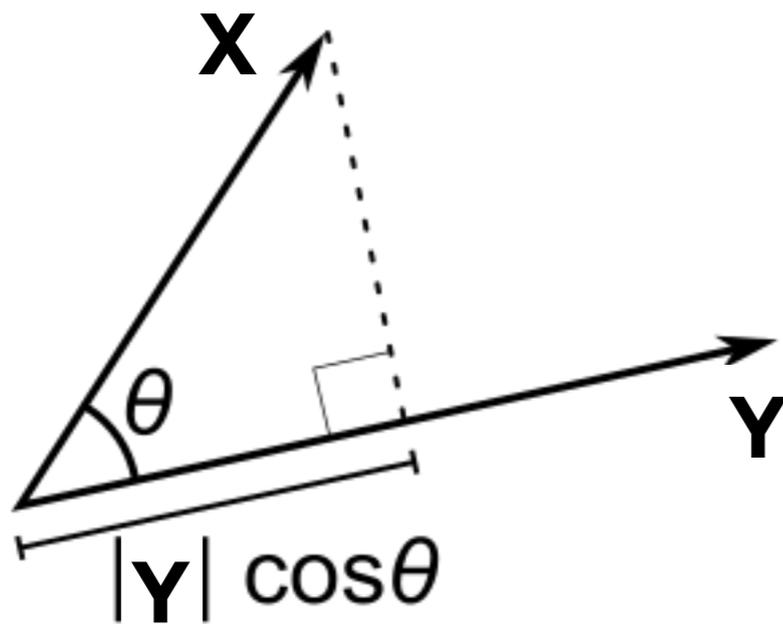
$$\frac{X \cdot Y}{|X| |Y|} = \cos(\theta)$$



High School Geometry

$x \cdot y$ “Distance” both vectors travel in the same direction

$\frac{X \cdot Y}{|X| |Y|}$ Angle between the two vectors



If I travel along X instead of Y, how far will I be from Y

Projection of X onto Y is linear



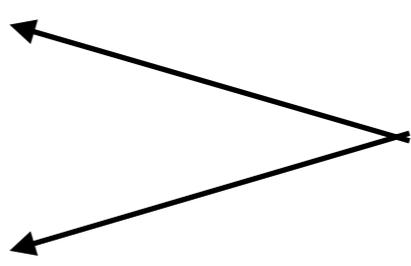
Back to Stats

X and Y are measured quantities of individuals in the same population.

$X = [1.5, 3.6, 2.3, \dots, 10.7]$

$Y = [2.9, 7.6, 4.3, \dots, 20.1]$

Can think of them as big vectors



If I travel along X instead of Y, how far will I be from Y

If I use X to “predict” Y, how close will I be



Make This More Precise

X and Y are measured quantities of individuals in the same population.

$$X = [1.5, 3.6, 2.3, \dots, 10.7]$$

$$X' = X - \text{mean}(X)$$

$$Y = [2.9, 7.6, 4.3, \dots, 20.1]$$

$$Y' = Y - \text{mean}(Y)$$

$$X' \cdot Y'$$

Covariance

$$\frac{X' \cdot Y'}{|X'| |Y'|}$$

Correlation

Exactly the
standard deviation!

$$S = \sqrt{\frac{\sum (X - \bar{X})^2}{N}}$$



Covariance and Correlation

The strength of “linear” dependence between two sets of observations.

$X' \cdot Y'$ Covariance

$$\text{Cov}(X, Y) = \frac{\sum (X_i - \bar{X})(Y_j - \bar{Y})}{n}$$

$\frac{X' \cdot Y'}{|X'| |Y'|}$ Correlation “Pearson”

$$\frac{\text{Cov}(x, y)}{\sigma_x * \sigma_y}$$



CHIDATA

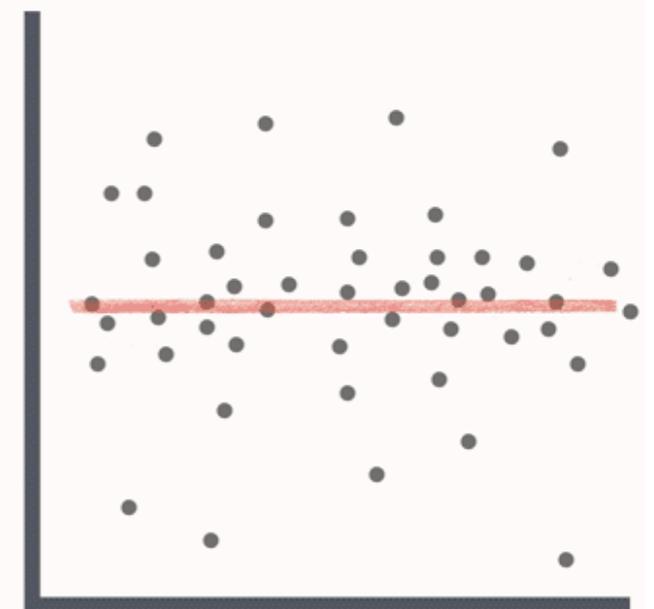
Covariance and Correlation



Positive Correlation



Negative Correlation



No Correlation



Recap: Significance

Control

Treatment

Measure how strongly correlated treatment is to a positive outcome.



Example

Hypothesis A proposed relationship between two variables.

$X = \{\text{treatment, control}\}$

$Y = \text{Weight Loss}$

$X = [0, 1, 0, \dots, 1, 1, 0]$

$Y = [1.7, -2.3, 7.6, \dots, 4.6]$

← **All of the data collected
in the trial**

$$\frac{X' \cdot Y'}{|X'| |Y'|} \sim 1 \quad \text{Strong association}$$

$$\frac{X' \cdot Y'}{|X'| |Y'|} \sim 0 \quad \text{Weak association}$$

$$\frac{X' \cdot Y'}{|X'| |Y'|} < 0 \quad \text{Counterproductive}$$



Extend to Random Variables

The strength of “linear” dependence between two random variables X and Y .

$$\begin{aligned} \text{cov}(X, Y) &= E[(X - E[X])(Y - E[Y])] \\ &= E[XY - XE[Y] - YE[X] - E[X][Y]] \\ &= E[XY] - E[X]E[Y] - E[Y]E[X] - E[X][Y] \\ &= E[XY] - E[X][Y] \end{aligned}$$



Properties

$$\text{Cov}(X, X) = \text{Var}(X)$$

$$\text{Cov}(X, \text{constant}) = 0$$

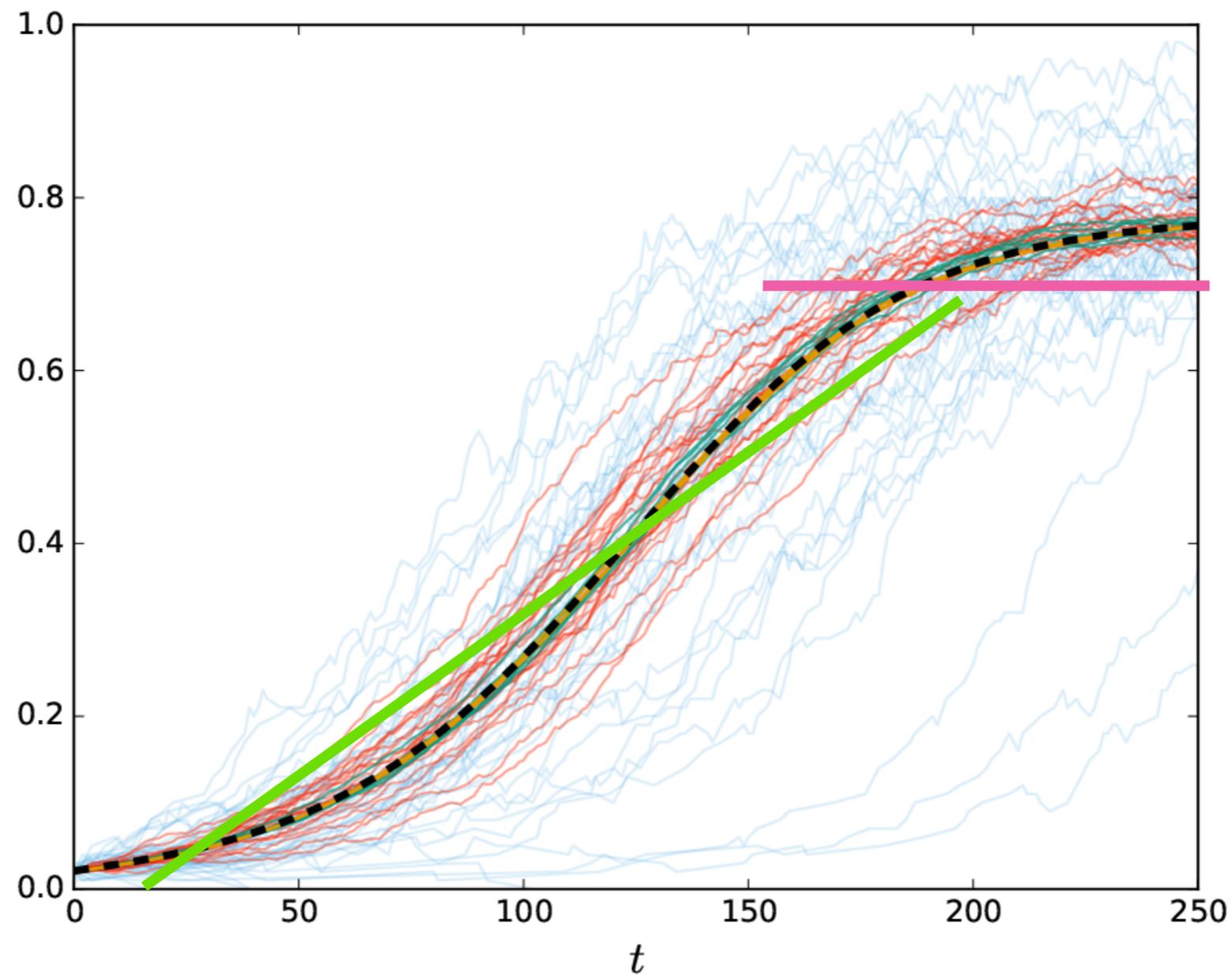
$$X, Y \text{ independent} \Rightarrow \text{Cov}(X, Y) = 0$$

!!!Converse not true!!!



Data Representation Matters

Saturation Effects





Data Representation Matters

Cities with > 10 Fire Stations

86 fires

101 fires

54 fires

Cities with ≤ 10 Fire Stations

3 fires

2 fires

4 fires

$X = \{>10 \text{ FSs}, \leq 10 \text{ FSs}\}$

$Y = \#\text{Fires}$



Data Representation Matters

Cities

86 fires, 15 FSs

3 fires, 1 FS

101 fires, 32 FSs

2 fires, 4 FSs

54 fires, 16 FSs

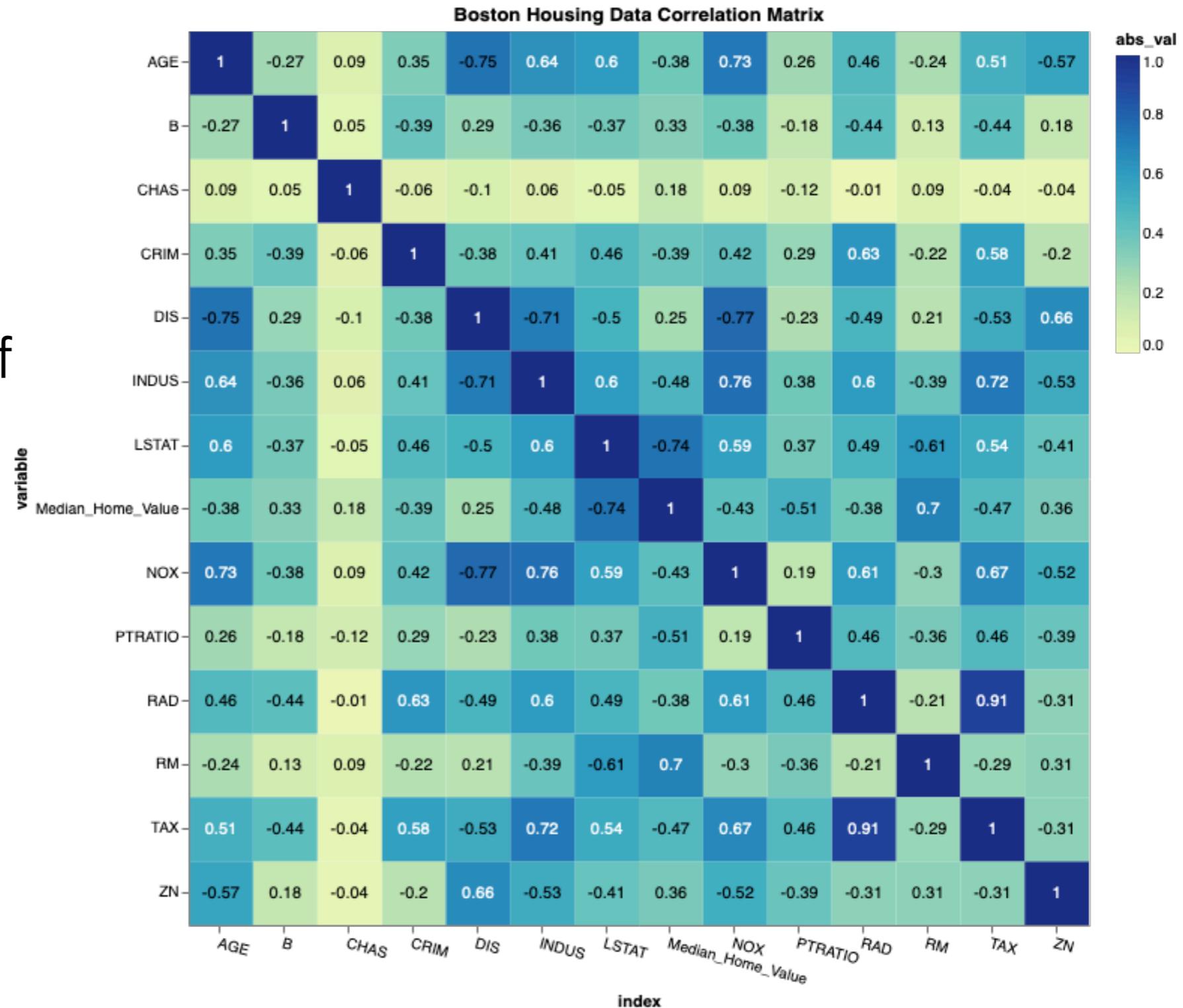
4 fires, 3 FSs

X = #FSs

Y = #Fires

Multi-Variable Correlation

Compare all pairs of variables





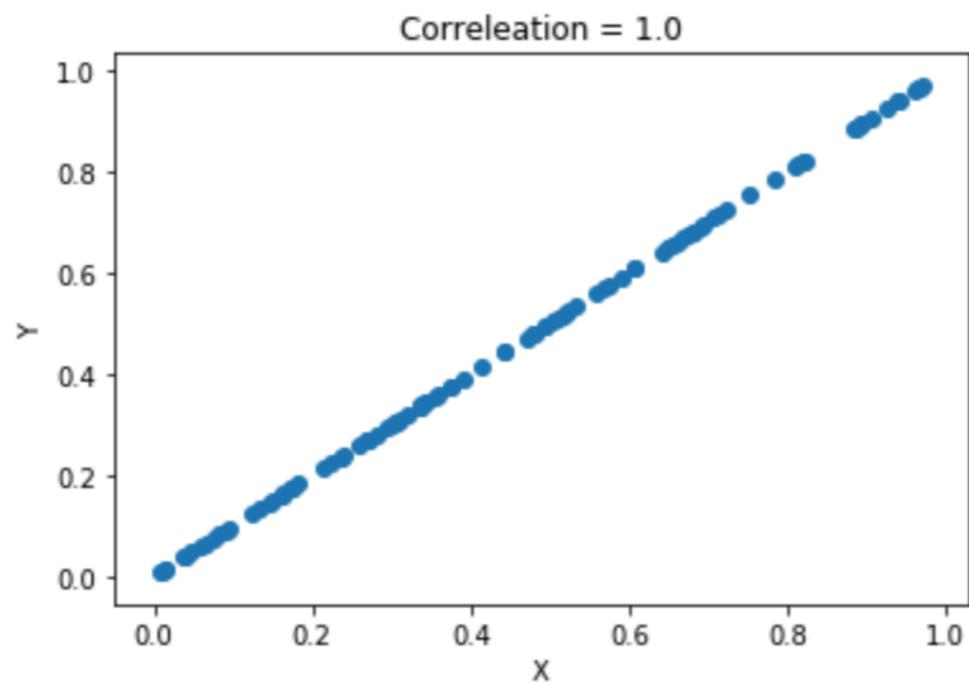
Examples

```
In [18]: %matplotlib inline
import numpy as np
import matplotlib.pyplot as plt
import scipy.stats

x = np.random.rand(100)
y = x

corr = scipy.stats.pearsonr(x,y)[0]

plt.scatter(x,y)
plt.xlabel('X')
plt.ylabel('Y')
plt.title('Correleation = '+ str(corr))
plt.show()
```

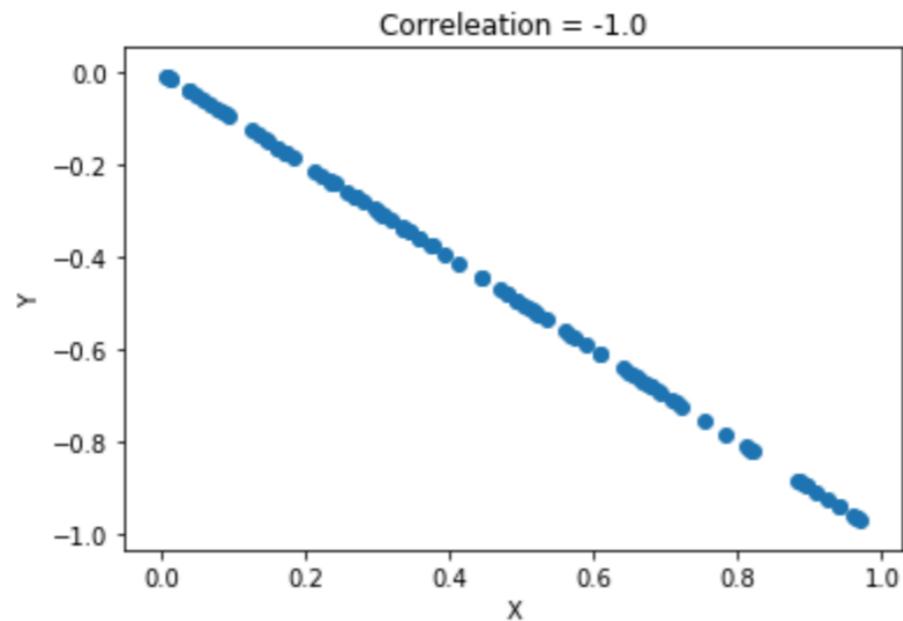




Examples

In [19]:

```
y = -x  
  
corr = scipy.stats.pearsonr(x,y)[0]  
  
plt.scatter(x,y)  
plt.xlabel('X')  
plt.ylabel('Y')  
plt.title('Correleation = '+ str(corr))  
plt.show()
```



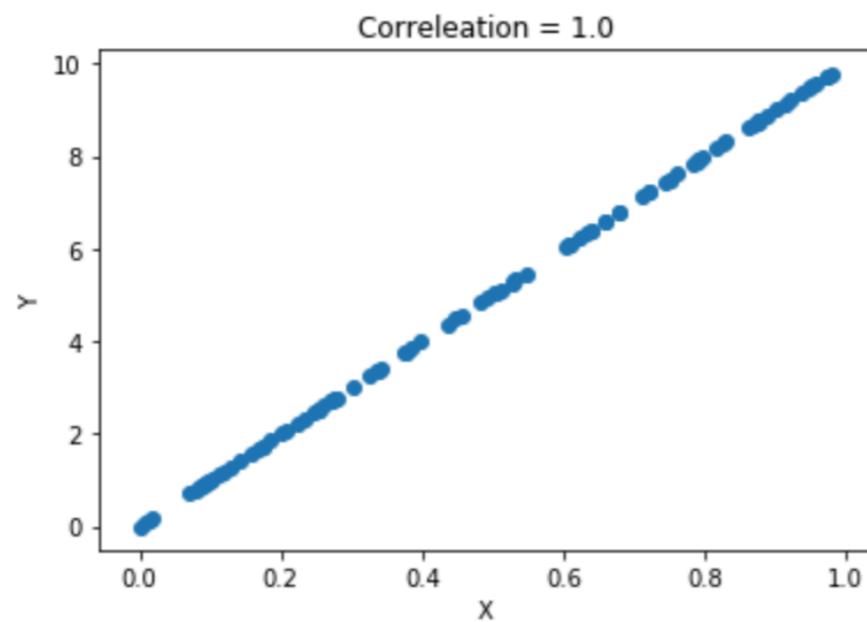


Examples

```
In [20]: x = np.random.rand(100)
y = 10*x
corr = scipy.stats.pearsonr(x,y)[0]

corr = scipy.stats.pearsonr(x,y)[0]

plt.scatter(x,y)
plt.xlabel('X')
plt.ylabel('Y')
plt.title('Correleation = '+ str(corr))
plt.show()
```



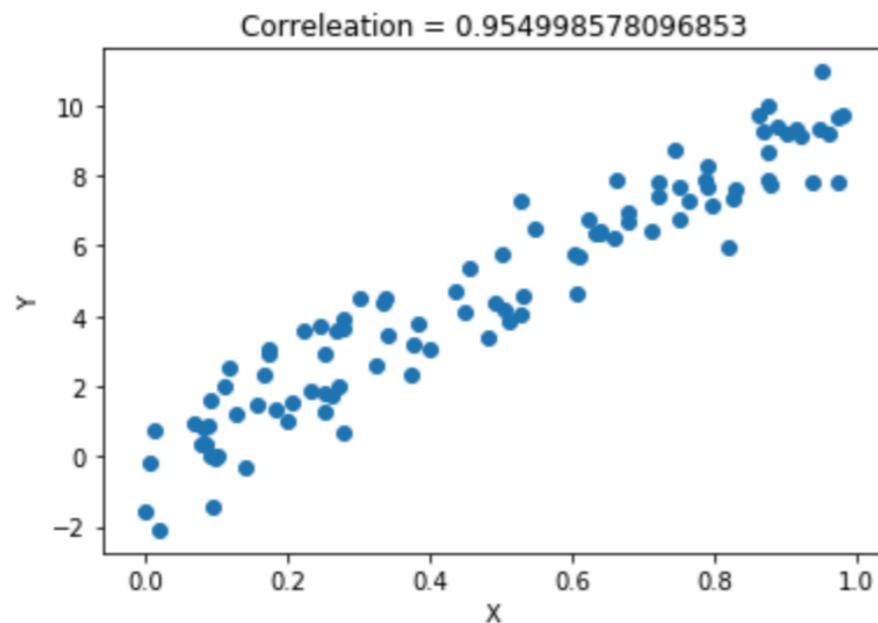


Examples

$$Y = Ax + \text{noise}$$

```
In [22]: y = 10*x + np.random.randn(100)
corr = scipy.stats.pearsonr(x,y)[0]
corr = scipy.stats.pearsonr(x,y)[0]

plt.scatter(x,y)
plt.xlabel('X')
plt.ylabel('Y')
plt.title('Correleation = ' + str(corr))
plt.show()
```



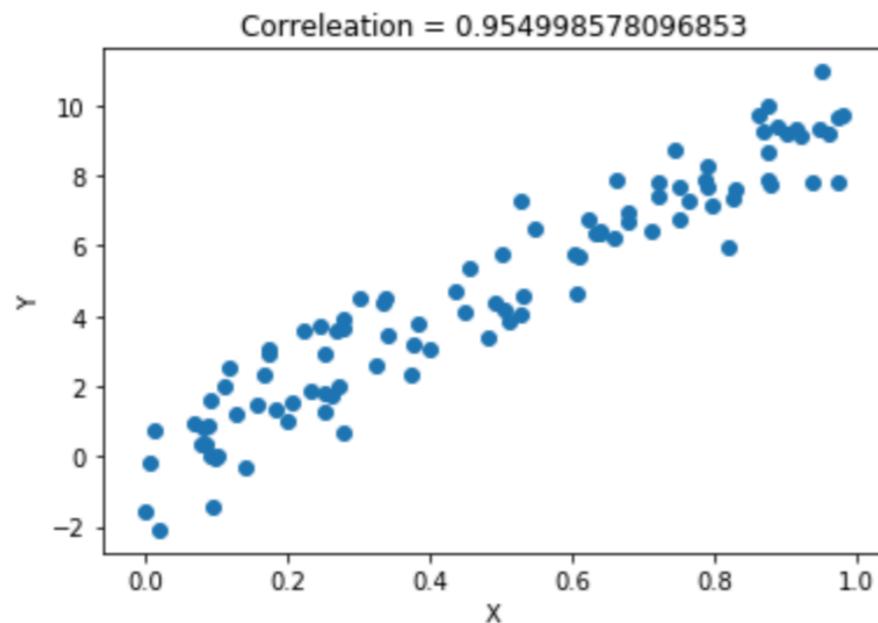


Examples

$$Y = Ax + \text{noise}$$

```
In [22]: y = 10*x + np.random.randn(100)
corr = scipy.stats.pearsonr(x,y)[0]
corr = scipy.stats.pearsonr(x,y)[0]

plt.scatter(x,y)
plt.xlabel('X')
plt.ylabel('Y')
plt.title('Correleation = ' + str(corr))
plt.show()
```





Examples

$$Y = A \cdot X + \epsilon \quad \epsilon \sim \text{i.i.d}$$

$$\text{Cov}(X, Y) = \mathbf{E}[XY] - \mathbf{E}[X]\mathbf{E}[Y]$$

$$\text{Cov}(X, Y) = \mathbf{E}[X(A X + \epsilon)] - \mathbf{E}[X]\mathbf{E}[(A X + \epsilon)]$$

$$\text{Cov}(X, Y) = \mathbf{E}[A X^2] + \mathbf{E}[X \epsilon] - \mathbf{E}[X]\mathbf{E}[A X] - \mathbf{E}[X]\mathbf{E}[\epsilon]$$

$$\underline{\text{Cov}(X, Y) = A \cdot \text{Var}[X]}$$



Examples

$$Y = A \cdot X + \epsilon \quad \epsilon \sim \text{i.i.d}$$

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}}$$

Exercise for you guys....

$$\text{Corr}(X, Y) = \frac{1}{\sqrt{1 + \frac{\text{Var}(\epsilon)}{A^2 \cdot \text{Var}(X)}}$$

$$\text{Corr}(X, Y) = \frac{1}{\sqrt{1 + \text{SNR}^{-1}}}$$

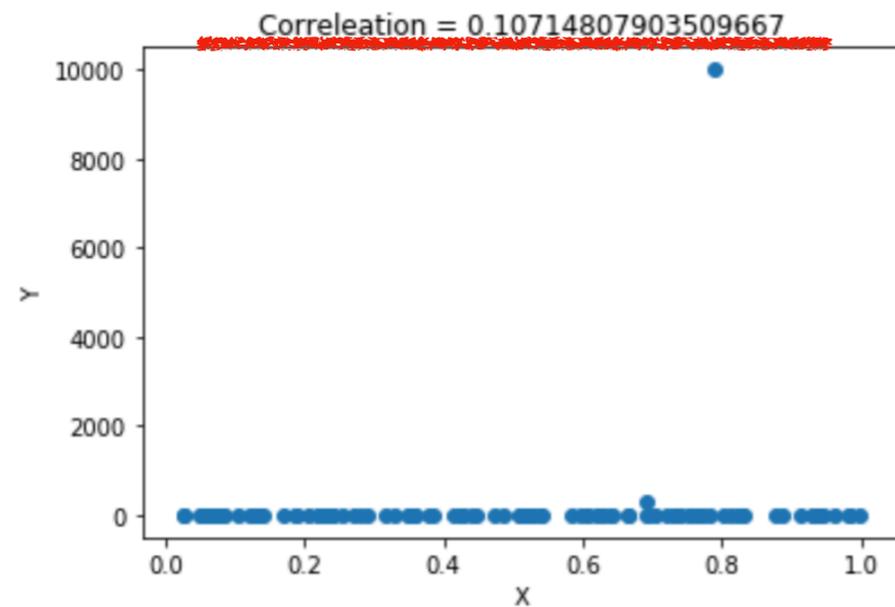


Examples

```
In [16]: %matplotlib inline
import numpy as np
import matplotlib.pyplot as plt
import scipy.stats

x = np.random.rand(100)
y = 10*x + np.random.randn(100)
y[0] = 300
y[1] = 10000
corr = scipy.stats.pearsonr(x,y)[0]

plt.scatter(x,y)
plt.xlabel('X')
plt.ylabel('Y')
plt.title('Correleation = ' + str(corr))
plt.show()
```





Summary: Covariance and Correlation

The strength of “linear” dependence between two sets of observations.

$X' \cdot Y'$ Covariance

$$\text{Cov}(X, Y) = \frac{\sum (X_i - \bar{X})(Y_j - \bar{Y})}{n}$$

$\frac{X' \cdot Y'}{|X'| |Y'|}$ Correlation “Pearson”

$$\frac{\text{Cov}(x, y)}{\sigma_x * \sigma_y}$$