

CS 2810 April 11

Admin:

- HW 8 due sunday (HW9 release monday)
- quiz4 scheduling (to be clear: you can still take it when sec3 scheduled, if you prefer it!)

Correlation

- what it is
- how its different than covariance

Causation

- its not correlation, even if its easily confused

Conditional Probability

- bayes rule
- independence

## Correlation (and covariance): Intuition (How two values vary together)

The behavior between any two values  $x$  and  $y$  can be summarized in one of the three ways:

1. as  $x$  gets larger  $y$  typically gets larger too

- ex:

- `x=temp on some day`

- `y=number of people on the beach on the same day`

- covariance & correlation is positive

2. as  $x$  gets larger  $y$  typically doesn't get larger or smaller

- ex:

- `x=individual's favorite number`

- `y=number of hot dogs that individual has eaten in their lifetime`

- covariance & correlation is zero

3. as  $x$  gets larger,  $y$  typically gets smaller

- ex:

- `x=average speed of driver on 10 mile commute`

- `y=average commute time of driver on 10 mile commute`

- covariance & correlation are negative

→ CORR MEASURES JUST THIS

→ COV INCLUDES OTHER FEATURES TOO

**So ... whats the difference between correlation and covariance?**

How does scaling X impact Sample Cov(X, Y)?

Let:

X be the hours spent doing a HW

Y be the score received on the HW

X	1	2	3
Y	70	80	90

$$\bar{X} = 2$$
$$\bar{Y} = 80$$

$$\hat{\sigma}_{xy}^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x})(y_i - \bar{y})$$

$$= \frac{1}{3-1} \left[ (1-2)(70-80) + (2-2)(80-80) + (3-2)(90-80) \right]$$

$$= \frac{1}{2} \left[ (-1)(-10) + (0)(0) + (1)(10) \right]$$
$$= 10$$

How does scaling X impact Sample Cov(X, Y)?

Let:

X be the minutes spent doing a HW

Y be the score received on the HW

X	60	120	180
Y	70	80	90

$$\bar{X} = 120$$

$$\bar{Y} = 80$$

$$\hat{\sigma}_{xy}^2 = \frac{1}{n-1} \sum_i (x_i - \bar{x})(y_i - \bar{y})$$

$$= \frac{1}{3-1} \left[ (60-120)(70-80) + (120-120)(80-80) + (180-120)(90-80) \right]$$

$$= \frac{1}{2} \left[ (-60)(-10) + (0)(0) + (60)(10) \right]$$

$$= 600$$

the magnitude of the  $\text{cov}(x, y)$  does not measure the strength of the correlation between  $x, y$

How does scaling X impact Sample Cov(X, Y)?

Let c be a scalar while  $(x_i, y_i)$  are paired observations.

$$\begin{aligned} \text{SAMPLE COV OF } cX_i \text{ AND } Y_i &= \frac{1}{N-1} \sum_i (cx_i - c\bar{x})(y_i - \bar{y}) \\ &= c \hat{\sigma}_{xy} \end{aligned}$$

SCALING (SAMPLES OR OUTCOMES) OF X BY C  
WILL SCALE COVARIANCE BY C TOO!

**!!Super big problem!!**

**Changing the units we measure data in will change our covariance scales**

- sample cov(hour study, grade) = 10
- sample cov(min study, grade) = 600
- ... we shouldn't interpret the magnitude of covariance as a "strength" of correlation!

**Covariance measures three things:**

- standard deviation of x
- standard deviation of y
- correlation between x and y (see "intuition" slide above)

**We want a "scale invariant" way of measuring correlation**

**Scale invariant - a statistic which doesn't change when data is scaled**



## A BOUND ON COVARIANCE

$$|\text{COV}(X, Y)| = |E[(X - E[X])(Y - E[Y])]|$$

$$\leq |E[(X - E[X])^2]^{1/2} E[(Y - E[Y])^2]^{1/2}|$$

$$= \sqrt{\text{VAR}(X) \text{VAR}(Y)} = \sigma_X \sigma_Y$$

COV MAY NOT BE GREATER IN MAGNITUDE THAN GEOMETRIC MEAN OF  $\text{VAR}(X)$   $\text{VAR}(Y)$

$$|\text{COV}(XY)| \leq \sigma_x \sigma_y$$

$$\text{SMALLEST } \text{COV}(XY) = -\sigma_x \sigma_y$$

$$\text{BIGGEST } \text{COV}(XY) = \sigma_x \sigma_y$$

# CORRELATION (AKA PEARSON'S CORRELATION COEFFICIENT)

CORRELATION

$$\rho_{xy} = \frac{\text{COV}(x, y)}{\sigma_x \sigma_y}$$

How MUCH COV DO X, Y HAVE COMPARED TO WHAT THEY COULD HAVE?

$$-1 \leq \rho_{xy} \leq 1$$

$+\sigma_x \sigma_y$  IS MAX COV(x, y)  
 $-\sigma_x \sigma_y$  IS MIN COV(x, y)

ESTIMATOR:

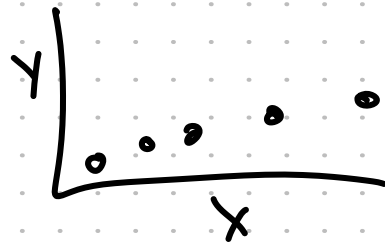
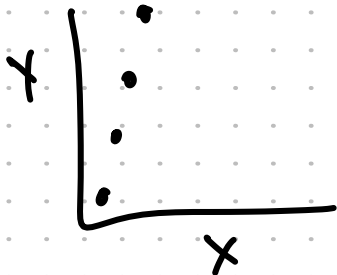
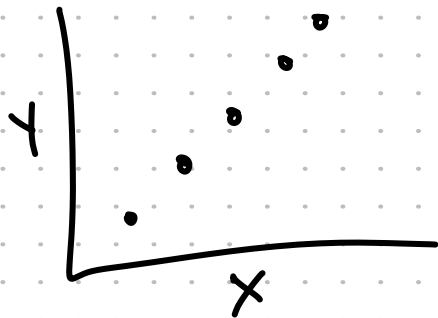
$$\hat{\rho}_{xy} = \frac{\hat{\sigma}_{xy}}{\hat{\sigma}_x \hat{\sigma}_y}$$

# MAX CORRELATION

$$\text{BOUND } |\text{cov}(x, y)| \leq \sigma_x \sigma_y$$

$$1 = \rho_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} \Rightarrow \text{cov}(x, y) = \sigma_x \sigma_y$$

COV IS AS BIG AS POSSIBLE GIVEN SCALE OF X AND Y



All examples have max correlation = 1. This implies that when x increases by some amount then y will always increase by m times that amount. (m is positive)

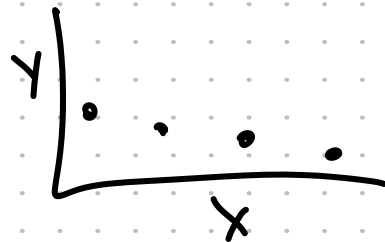
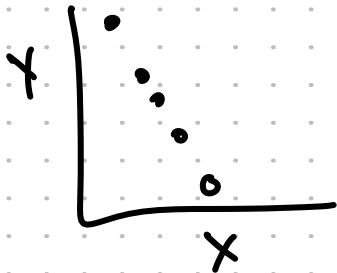
# MIN CORRELATION

$$\text{BOUND } |\text{cov}(x, y)| \leq \sigma_x \sigma_y$$

$$-1 = \rho_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} \Rightarrow$$

$$\text{cov}(x, y) = -\sigma_x \sigma_y$$

COV IS AS SMALL AS POSSIBLE GIVEN SCALE OF X AND Y



All examples have min correlation = -1. This implies that when x increases by some amount then y will always decrease by m times that amount (m is negative)

ICA 1

$$\hat{\sigma}_y^2 = \frac{1}{3} \left[ (2-2)^2 + (2-2)^2 + (1-2)^2 + (3-2)^2 \right]$$

= 2/3

Compute the sample correlation for the four samples below.

Give one sentence which interprets its meaning so a non-technical reader can easily understand.

X	7	5	5	3
Y	2	2	1	3

$$\bar{x} = 5 \quad \bar{y} = 2$$

$$\hat{\sigma}_x^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$$

$$= \frac{1}{4-1} \left[ (7-5)^2 + (5-5)^2 + (5-5)^2 + (3-5)^2 \right] = \frac{1}{3} (4+0+0+4) = \frac{8}{3}$$

$$\hat{\sigma}_{xy} = \frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y})$$

$$= \frac{1}{4-1} \left[ (7-5)(2-2) + (5-5)(2-2) + (5-5)(1-2) + (3-5)(3-2) \right]$$

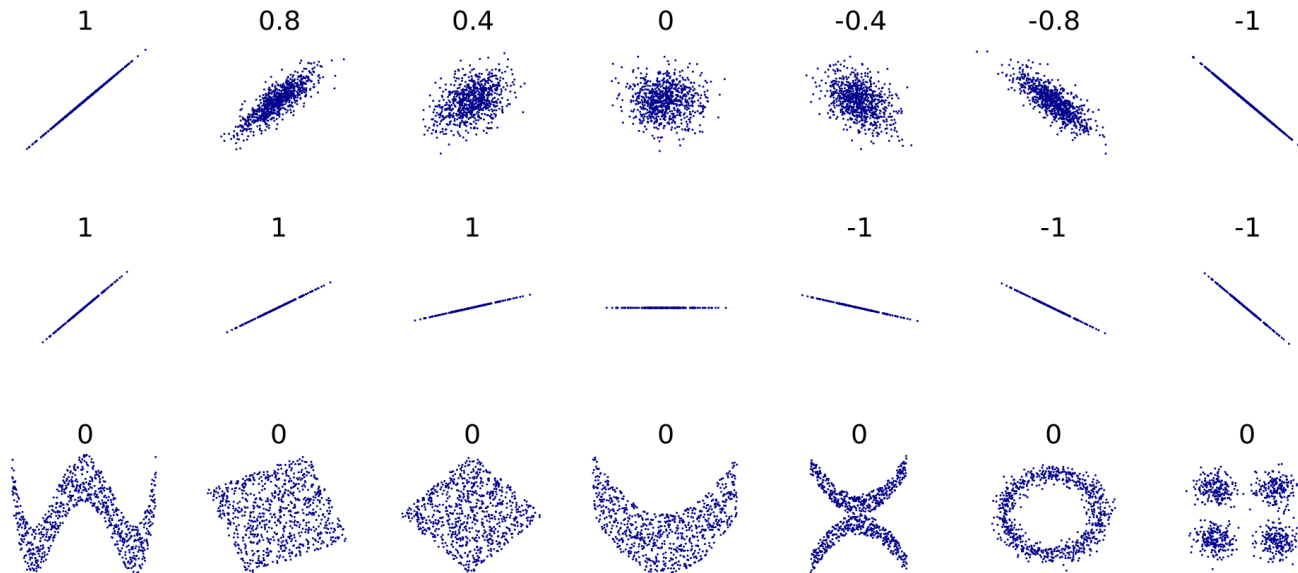
$$= \frac{1}{3} [0+0+0-2] = -\frac{2}{3}$$

$$\hat{\rho}_{xy} = \frac{\hat{\sigma}_{xy}}{\hat{\sigma}_x \hat{\sigma}_y} = \frac{-2/3}{\sqrt{8/3} \sqrt{2/3}} = -\frac{2}{2\sqrt{2}} = -\frac{1}{\sqrt{2}}$$

## What does correlation mean (and what it doesn't) (consistency)

Correlation describes the consistency with which one variable changes in response to another.

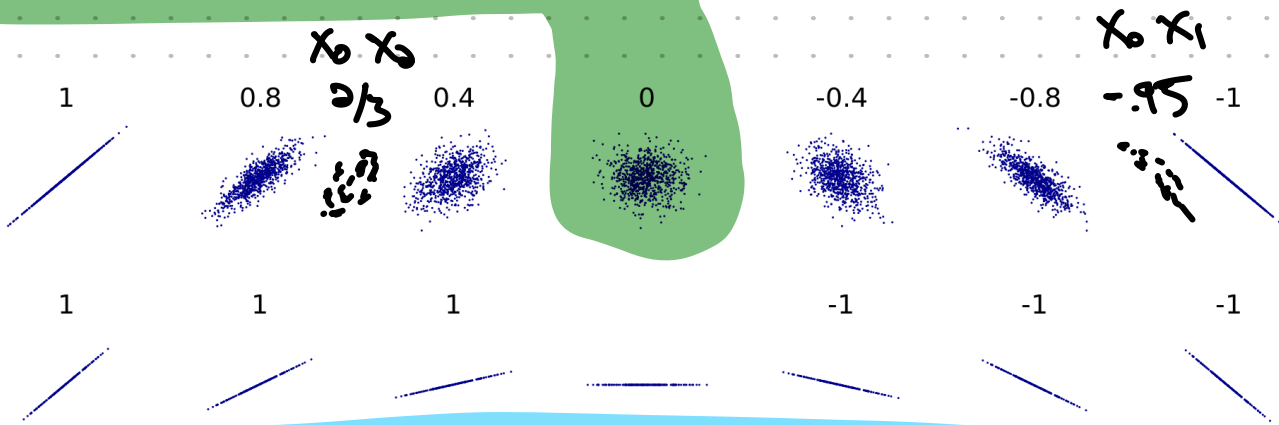
The most consistent relationship means whenever one variable goes up by one, the other always changes by the same amount (forming a line of observations).



CREDIT  
WIKIPEDIA

# What does correlation mean (and what it doesn't) (independence)

If  $x$  and  $y$  are independent then  $\text{corr}(x, y) = 0$



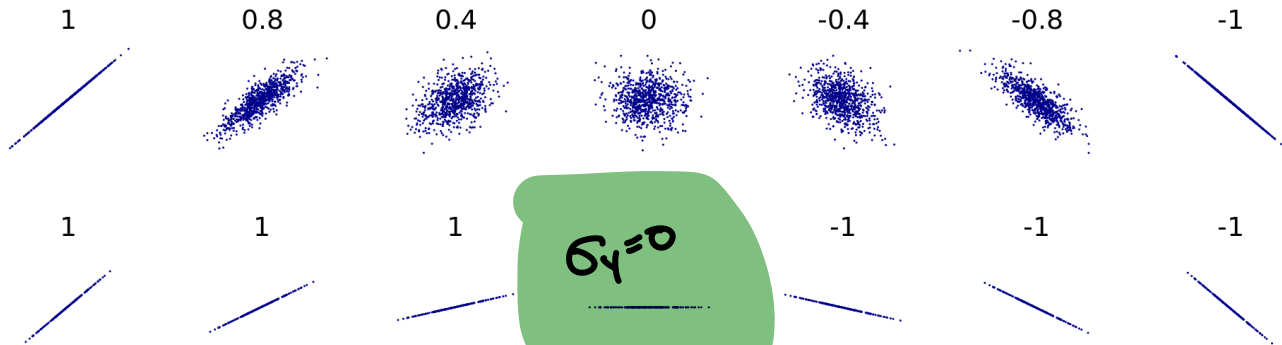
If  $\text{corr}(x, y) = 0$ , it may not be true that  $x$  and  $y$  are independent (see all examples below)





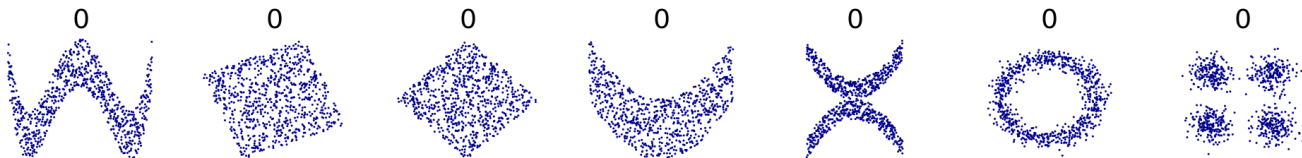
What does correlation mean (and what it doesn't) (0 variance case)

Correlation is only defined when each  $x, y$  has a positive variance (we'd divide by zero otherwise)



$$\rho_{xy} = \frac{\text{cov}(xy)}{\sigma_x \sigma_y}$$

If  $\text{corr}(x, y) = 0$ , it may not be true that  $x$  and  $y$  are independent (see all examples below)



**"strength" of a correlation is the distance from zero (in either direction)**

# CLASS PRACTICE

For each pair of random variables below, estimate Pearson's correlation coefficient.

$x$  = # hours student studies for final,  $y$  = final grade of that student.

$x$  = favorite number of student,  $y$  = student's best friend's favorite number.

$x$  = how many quarters one spends at supermarket,  $y$  = weight after - weight before.

$x$  = miles run in the past year,  $y$  = average speed of a runner on a 1 mile course.

$x$  = internet explorer usage rate,  $y$  = covid cases in US (per year ... past 30 years).

## In Class Assignment 2

Compute correlation between  $X_0, X_1$

$$\rho_{X_0 X_1} = \frac{\text{COV}(X_0, X_1)}{\sigma_{X_0} \sigma_{X_1}} = \frac{-1.9}{\sqrt{1} \sqrt{4}} = -0.95$$

Which feature shows the most consistent relationship with  $X_0$ :  $X_1$  or  $X_2$ ?

$$\Sigma = \begin{bmatrix} 1 & -1.9 & 2 \\ -1.9 & 4 & \text{---} \\ 2 & \text{---} & 9 \end{bmatrix}$$

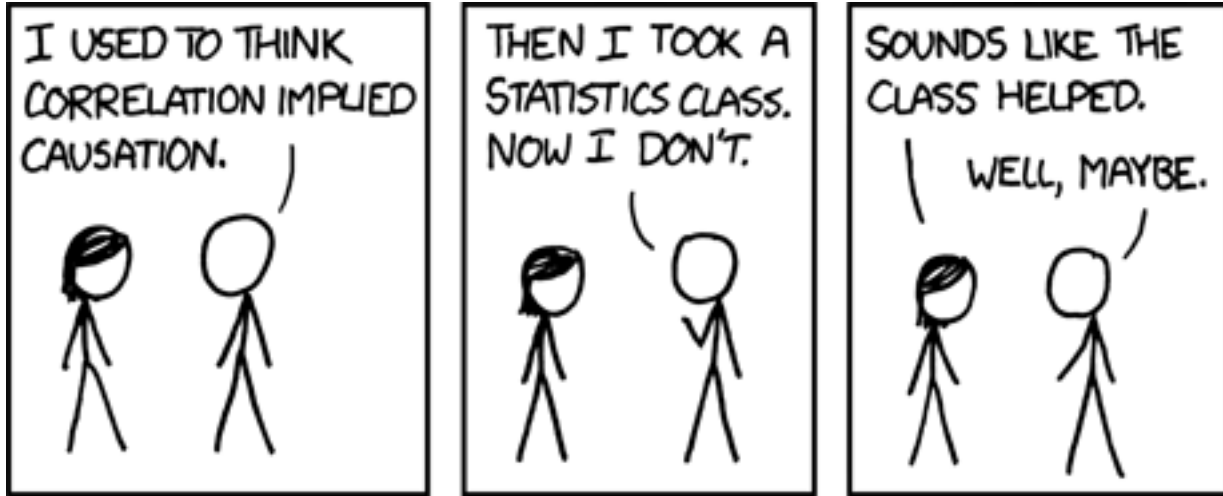
NOT GIVEN

$$= \begin{bmatrix} \text{COV}(X_0, X_0) & \text{COV}(X_0, X_1) & \text{COV}(X_0, X_2) \\ \text{COV}(X_1, X_0) & \text{COV}(X_1, X_1) & \text{COV}(X_1, X_2) \\ \text{COV}(X_2, X_0) & \text{COV}(X_2, X_1) & \text{COV}(X_2, X_2) \end{bmatrix}$$

$$\text{COV}(X_0, X_0) = \text{VAR}(X_0) = \sigma_{X_0}^2$$

$$\rho_{X_0 X_2} = \frac{2}{\sqrt{1} \sqrt{9}} = \frac{2}{3}$$

## Correlation does not imply causality



CREDIT:  
XKCD

see silly examples here: <https://www.tylervigen.com/spurious-correlations>

all these correlations are indeed non-zero and show some relationship between variables, though we shouldn't expect that changing  $x$  necessarily impact  $y$   
(discuss titanic hw problem)

## Bayes Rule & Conditional Probability

## Conditional Probability (intuition / definition)

Whats the probability that a person has covid?

$$P(C=1)$$

Whats the probability that a person has covid given a positive antigen test?

$$P(C=1 | A=1)$$

Whats the probability that a person has covid given a negative antigen test?

$$P(C=1 | A=0)$$

Whats the probability that an antigen test is negative given a person has covid?

$$P(A=0 | C=1)$$

Whats the probability that an antigen test is positive given a person doesn't have covid?

$$P(A=1 | C=0)$$

A conditional probability gives the probability of one event given that another has occurred.

Let C be a random variable representing whether a person has covid (1=covid, 0=no covid)

Let A be a random variable representing whether that person's antigen test is positive (1=positive, 0=negative)

Conditional Probability (algebraic definition)

PROB A HAPPENS  
GIVEN B ALREADY  
WAS

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

PROB A, B HAPPEN  
TOGETHER

PROB B HAPPENS

A = 1 IF STUDENT HAS  
WHITE SOCKS  
0 OTHERWISE  
B = 1 STUDENT ON PROF RIGHT  
SIDE OF CLASS

$$P(A=1|B=1) = \frac{P(A=1 \cap B=1)}{P(B=1)} = \frac{17}{30}$$





### In Class Assignment 3


Estimate each of the probabilities to the right from the observation table below. Give a sentence which explains their meaning so a non-technical reader could easily understand their value.

### BITCOIN PRICE MOVEMENT

TWITTER SENTIMENT SCORE  
(TWEETS w/ BITCOIN HASHTAG)

$S=1$  

$S=0$  

$S=-1$  

	$B=1$ UP	$B=-1$ DOWN
$S=1$	10K	9K
$S=0$	4K	7K
$S=-1$	2K	11K


①  $P(S=-1 \mid B=1) = \frac{2}{10+9+4+7+2+11}$   
 $= 2/43$


### In Class Assignment 3


Estimate each of the probabilities to the right from the observation table below. Give a sentence which explains their meaning so a non-technical reader could easily understand their value.

### BITCOIN PRICE MOVEMENT

TWITTER SENTIMENT SCORE (TWEETS w/ BITCOIN HASHTAG)

$S=1$  

$S=0$  

$S=-1$  

	$B=1$ UP	$B=-1$ DOWN
$S=1$	10K	9K
$S=0$	4K	7K
$S=-1$	2K	11K

$$\textcircled{2} P(S=-1 | B=1) = \frac{2}{10+4+2} = \frac{2}{16}$$

### In Class Assignment 3

Estimate each of the probabilities to the right from the observation table below. Give a sentence which explains their meaning so a non-technical reader could easily understand their value.

BITCOIN PRICE MOVEMENT

	B=1 UP	B=-1 DOWN
S=1 😊	10K	9K
S=0 😐	4K	7K
S=-1 😞	2K	11K

TWITTER SENTIMENT SCORE (TWEETS w/ BITCOIN HASHTAG)

③  $P(B=1 | S=-1) = \frac{2}{2+11} = \frac{2}{13}$

④  $P(B=1) = \frac{10+4+2}{43} = \frac{16}{43}$

Now ARE ① ② ④ RELATED?

### In Class Assignment 3

Estimate each of the probabilities to the right from the observation table below. Give a sentence which explains their meaning so a non-technical reader could easily understand their value.

$$P(S=-1 | B=1) = \frac{P(S=-1, B=1)}{P(B=1)}$$

$$\frac{2}{16} = \frac{2/43}{16/43}$$

①  $P(S=-1, B=1)$

②  $P(S=-1 | B=1)$

④  $P(B=1)$

Now ARE ① ② ④ RELATED?