

# Midterm 2

- This Friday, 12-12:50pm  
MS4000A  
(here).

## Math 170E: Winter 2023

- Content: Focus on  
lectures

Lecture 19, Mon 27th Feb

$$10 \leq X \leq 18$$

↑  
Last Friday

The correlation coefficient

- 2 sided Cheat sheet,  
Simple Calculator.
- Practice on Canvas
- No lecture Friday
- My OH: Wed 5:30-6:30  
(this week). Thu 9:30-5:30

## Last time:

- Let  $X, Y$  be a pair of discrete random variables taking values in sets  $S_X, S_Y \subseteq \mathbb{R}$
- Let  $S = S_X \times S_Y$  and let  $X, Y$  have joint PMF  $p_{X,Y}(x, y)$
- If  $g : S \rightarrow \mathbb{R}$ , we define

$$\mathbb{E}[g(X, Y)] = \sum_{(x,y) \in S} g(x, y) p_{X,Y}(x, y)$$

$\rightarrow p_{X,Y}(x,y) = p_X(x)p_Y(y)$

- If  $X, Y$  are independent and  $g : S_X \rightarrow \mathbb{R}$ ,  $h : S_Y \rightarrow \mathbb{R}$ , then

$$\mathbb{E}[g(X)h(Y)] = \mathbb{E}[g(X)]\mathbb{E}[h(Y)]$$

$\rightarrow \text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + 2\{\mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]\}$

- We have the Cauchy-Schwarz inequality:

$$|\mathbb{E}[XY]| \leq \sqrt{\mathbb{E}[X^2]\mathbb{E}[Y^2]}$$

# Today:

We'll discuss today:

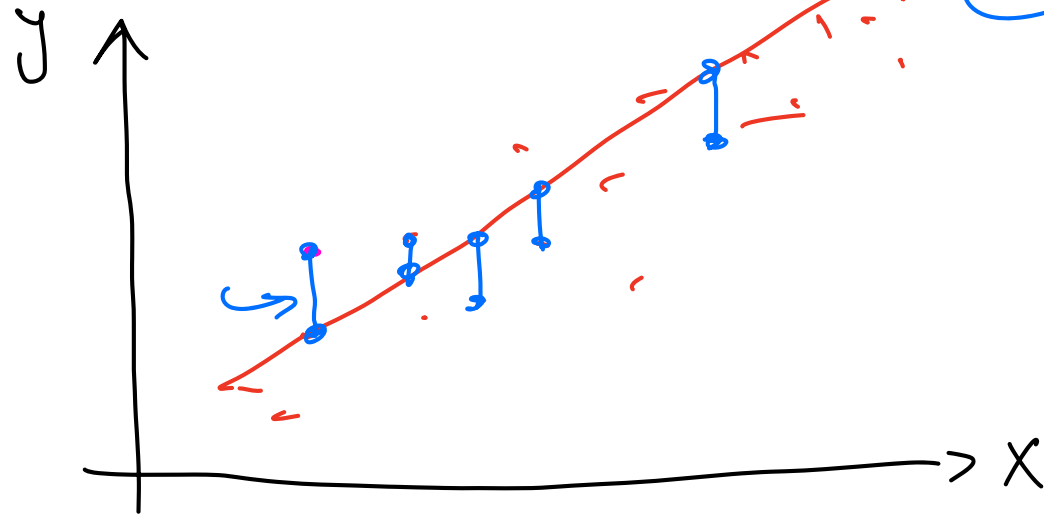
- the definition of the *covariance* and *correlation coefficient* of two random variables
- how to compute the *covariance* and *correlation coefficient* of two random variables
- the 'least-squares' line of best fit



Following problem:

Do an experiment & measured two variables  $(x, y)$

↳ Gives data points  $(x_j, y_j)$ ,  $j = 1, \dots, N$ .



↳ Plot indicates that there is a relationship b/w  $x$  &  $y$ . The simplest relationship is a linear one.

$$y = ax + b$$

for constants  $a, b$ . ↗ ↖ find some "good"  $a, b$ .

Ask for a minimal choice of  $(a, b)$ .

↳ We ask for  $(a, b)$  such that they minimise the sum of the vertical distances.

If a pair  $(a_N, b_N)$  minimises the above distance then we say that the line

$$y = a_N x + b_N$$

is a line of best fit.

How to find  $(a_N, b_N)$ ?

Minimise over  $(a, b) \in \mathbb{R}_N^2$ , the function:

$$d(a, b) = \sum_{i=1}^N (y_i - ax_i - b)^2$$

$$\hookrightarrow \nabla d(a, b) = (0, 0).$$

Solve for  $(a, b)$ :

$$\left\{ \begin{aligned} a_N &= \frac{\frac{1}{N} \sum_{i=1}^N x_i y_i - \left( \frac{1}{N} \sum_{i=1}^N x_i \right) \left( \frac{1}{N} \sum_{i=1}^N y_i \right)}{\frac{1}{N} \sum_{i=1}^N x_i^2 - \left( \frac{1}{N} \sum_{i=1}^N x_i \right)^2} \\ b_N &= \frac{1}{N} \sum_{i=1}^N y_i - \left( \frac{1}{N} \sum_{i=1}^N x_i \right) a_N \end{aligned} \right.$$

→ by Cauchy-Schwarz inequality  $> 0$

So sign of  $a_N$  purely depends on the sign of

$$\frac{1}{N} \sum_{i=1}^N x_i y_i - \left( \frac{1}{N} \sum_{i=1}^N x_i \right) \left( \frac{1}{N} \sum_{i=1}^N y_i \right)$$

If we think of  $X$ 's generated by a r.v.  $X$   
 $Y$ 's —————  $Y$

then

as  $N \rightarrow +\infty$ , should converge to:

$$\#[XY] - \#[X]\#[Y]$$

(Law of large  
Numbers  
Statement).

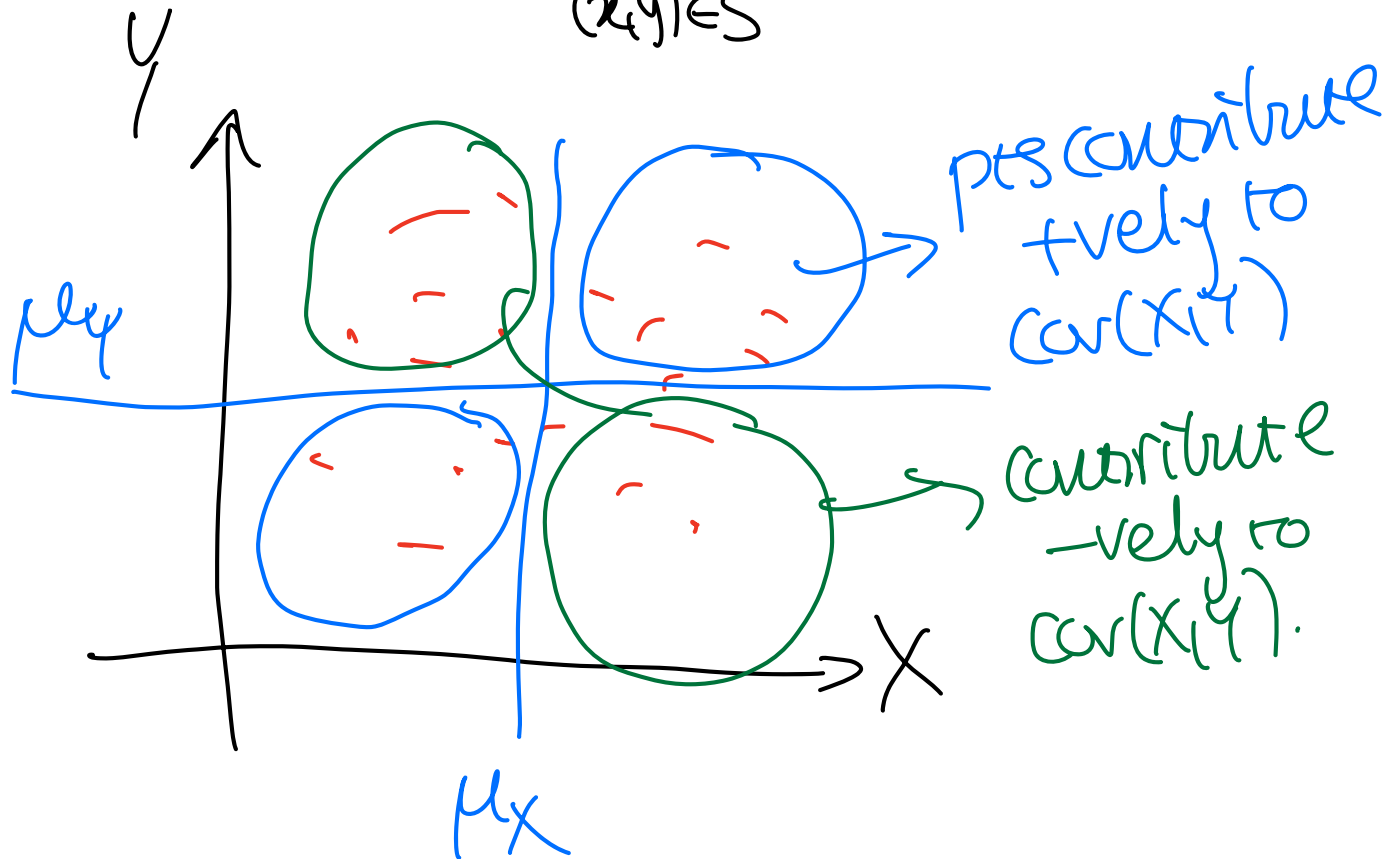
$$\left| \sum_{n=1}^N a_n b_n \right| \leq \left( \sum_{n=1}^N a_n^2 \right)^{1/2} \left( \sum_{n=1}^N b_n^2 \right)^{1/2} \dots$$

**Definition 4.7:** Let  $X, Y$  be a pair of discrete random variables taking values.

We define the **covariance** of  $X, Y$  to be

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

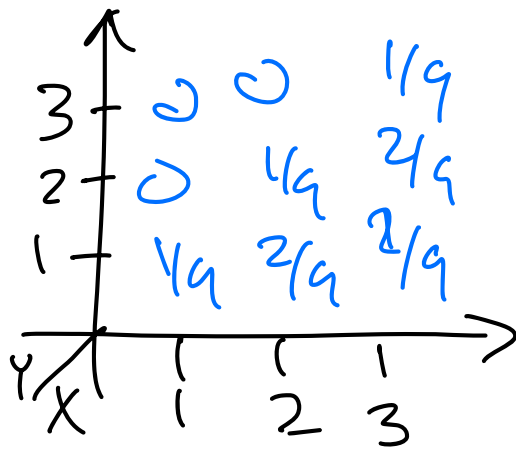
$$= \sum_{(x,y) \in S} (x - \mu_X)(y - \mu_Y) P_{X,Y}(x,y)$$





### Example 7:

- You choose two numbers at random from the set  $\{1, 2, 3\}$
- Let  $X$  be the larger and  $Y$  be the smaller of these two numbers
- What is  $\text{cov}(X, Y)$ ?



Last time

$$\hookrightarrow E[XY] = 4$$
$$E[X] = \frac{22}{9}$$
$$E[Y] = \frac{14}{9}$$

$$\hookrightarrow \text{cov}(X, Y) = 4 - \frac{22}{9} \times \frac{14}{9} = \frac{16}{81}.$$

**Proposition 4.8:** If  $X$  is a random variable, then

$$\text{cov}(X, X) = \text{var}(X).$$

**Proof:**

$$\begin{aligned}\text{cov}(X, X) &= E[(X - \mu_X)(X - \mu_X)] \\ &= E[(X - \mu_X)^2] = \text{var}(X).\end{aligned}$$

**Proposition 4.9:** If  $X, Y$  are *independent*, then

$$\text{cov}(X, Y) = 0$$

**Proof:**

$$\text{If } X, Y \text{ indep, } E(XY) = \mu_X \mu_Y.$$

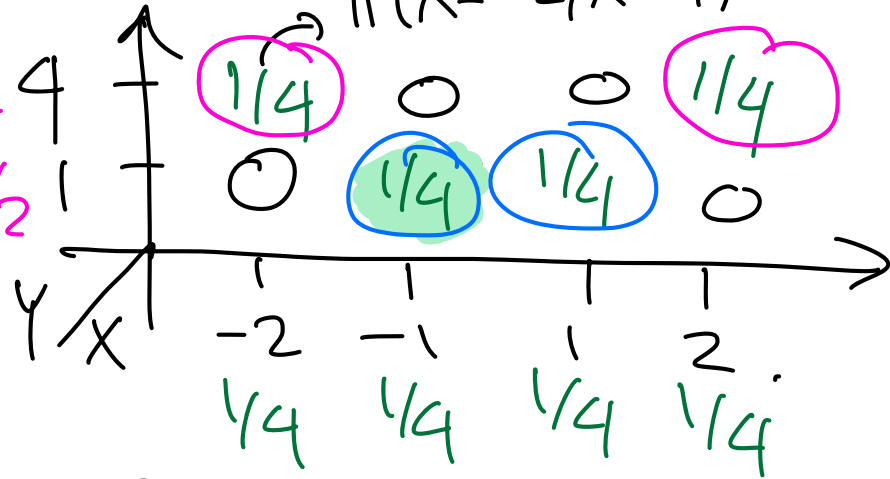
$$\text{If } \underline{\text{indep}} \Rightarrow \text{cov}(X, Y) = 0$$



### Example 8:

- Let  $X$  be a discrete r.v. which is uniform on  $\{-2, -1, 1, 2\}$  and set  $Y = X^2$
- What is  $\text{cov}(X, Y)$ ?
- Are  $X$  and  $Y$  independent?

$$P_{X,Y}(x,y) = P(X=x, X^2=y)$$



Now

$$E[XY] = \sum xy P_{X,Y}(x,y) = 0$$

$$E[X] = 0 \Rightarrow \text{cov}(X,Y) = 0.$$

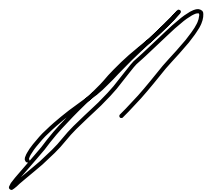
But  $X$  &  $Y$  are not independent!

$$P_{X,Y}(-1,1) = 1/4 \neq \frac{1}{2} \times \frac{1}{4} = P_X(-1)P_Y(1)$$

**Proposition 4.10:** Let  $X, Y$  be (discrete) random variables and  $a, b \in \mathbb{R}$ . Then,

$$\text{cov}(aX, bY) = ab \text{cov}(X, Y)$$

**Proof:**

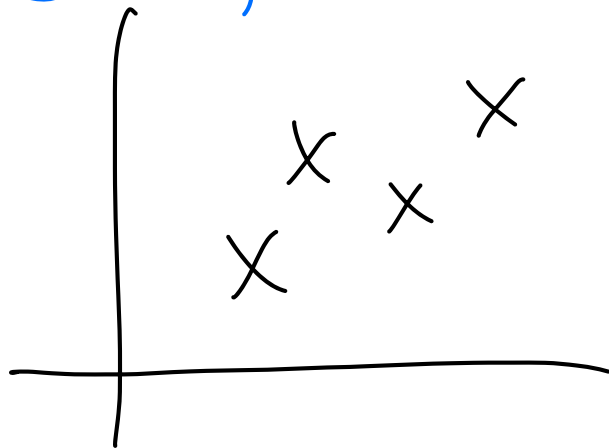
$$\begin{aligned} \text{cov}(aX, bY) &= E[abXY] - E[aX]E[bY] \\ &= ab \text{cov}(X, Y). \end{aligned}$$


## Definition 4.11:

- Let  $X, Y$  be a pair of (discrete) random variables
- We define the **correlation coefficient** of  $X, Y$  to be

(Pearson)

$$\rho(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)\text{var}(Y)}}$$

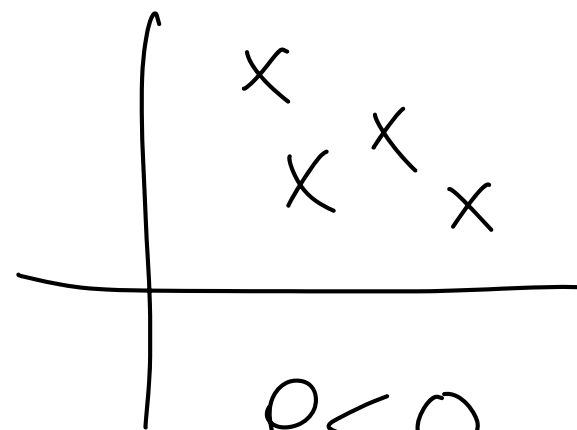


$\rho > 0$

Positively correlated

$$\rho(aX, bY) = \frac{\text{cov}(aX, bY)}{\sqrt{\text{var}(aX)\text{var}(bY)}}$$

$$\text{var}(aX) = a^2 \text{var}(X)$$



$\rho < 0$

Negatively correlated.

$$\begin{aligned} &= \frac{ab \text{cov}(X, Y)}{\sqrt{a^2 b^2} \sqrt{\text{var}(X)\text{var}(Y)}} = \frac{ab}{|a| |b|} \rho(X, Y) \\ &= \frac{ab}{|ab|} \rho(X, Y) \end{aligned}$$

**Proposition 4.12:**

If  $X, Y$  are (discrete) random variables, then

$\left[ \begin{array}{l} \text{if } X, Y \text{ are indep,} \\ \text{then } \rho(X, Y) = 0 \end{array} \right]$

$$-1 \leq \rho(X, Y) \leq 1$$

In particular,  $|\rho(X, Y)| = 1$  **if and only if** then  $Y = aX + b$  for some  $a, b \in \mathbb{R}$ .

**Proof:**

$$|\text{cov}(X, Y)| = \left| \mathbb{E}[(\overbrace{X - \mu_X}^W)(\overbrace{Y - \mu_Y}^Z)] \right|$$

$$\begin{aligned} \text{Cauchy-Schwarz} &\leq \sqrt{\mathbb{E}[(X - \mu_X)^2] \mathbb{E}[(Y - \mu_Y)^2]} \\ &= \sqrt{\text{Var}(X) \text{Var}(Y)} \end{aligned}$$

$$|\rho(X, Y)| = \frac{|\text{cov}(X, Y)|}{\sqrt{\text{Var}(X) \text{Var}(Y)}} \leq 1.$$

