

**Question 1**

- (a) Prove that for any random variable  $X$  and any constant  $a \in \mathbb{R}$ ,

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

- (b) Prove that for any random variable  $X, Y$  and  $Z$ , and constants  $a, b \in \mathbb{R}$ ,

$$\text{Cov}(aX + bY, Z) = a\text{Cov}(X, Z) + b\text{Cov}(Y, Z).$$

- (c) For any random variables  $X$  and  $Y$ , and constants  $a, b \in \mathbb{R}$ , find an expression for

$$\text{Cov}(aX + b, Y)$$

in terms of  $\text{Cov}(X, Y)$ .

**Solution to Question 1****Part (a):**

Recall the definition of covariance:

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

where  $\mu_X = \text{E}[X]$  and  $\mu_Y = \text{E}[Y]$ . Since  $a$  is a constant,  $\text{E}[a] = a$ . Therefore,

$$\text{Cov}(X, a) = \text{E}[(X - \mu_X)(a - a)] = \text{E}[0] = 0.$$

**Part (b):**

Using the linearity of expectation,

$$\mu_{aX+bY} = \text{E}[aX + bY] = a\text{E}[X] + b\text{E}[Y] = a\mu_X + b\mu_Y$$

Therefore, writing  $\mu_Z = \text{E}[Z]$ , and using the linearity of expectation

$$\begin{aligned} \text{Cov}(aX + bY, Z) &= \text{E}[(aX + bY - \mu_{aX+bY})(Z - \mu_Z)] \\ &= \text{E}[(aX + bY - a\mu_X - b\mu_Y)(Z - \mu_Z)] \\ &= \text{E}[(aX - a\mu_X + bY - b\mu_Y)(Z - \mu_Z)] \\ &= \text{E}[(aX - a\mu_X)(Z - \mu_Z) + (bY - b\mu_Y)(Z - \mu_Z)] \\ &= \text{E}[(aX - a\mu_X)(Z - \mu_Z)] + \text{E}[(bY - b\mu_Y)(Z - \mu_Z)] \\ &= \text{E}[a(X - \mu_X)(Z - \mu_Z)] + \text{E}[b(Y - \mu_Y)(Z - \mu_Z)] \\ &= a\text{E}[(X - \mu_X)(Z - \mu_Z)] + b\text{E}[(Y - \mu_Y)(Z - \mu_Z)] \\ &= a\text{Cov}(X, Z) + b\text{Cov}(Y, Z) \end{aligned}$$

**Part (c):**

Using Part (b),

$$\text{Cov}(aX + b, Y) = a\text{Cov}(X, Y) + \text{Cov}(b, Y)$$

and using Part (a),  $\text{Cov}(b, Y) = 0$ , which implies

$$\text{Cov}(aX + b, Y) = a\text{Cov}(X, Y).$$

(To justify the first part, you could consider  $b = bZ$ , where  $Z$  is the random variable which is identically equal to 1, and then proceed from there.)

## Question 2

Suppose you are tracking the value of two companies listed on the London Stock Exchange over the course of one week. Rather than record the actual values of the share prices, you record the increase or decrease in each share price value at the daily close to the nearest pound. You record the following table:

	Monday	Tuesday	Wednesday	Thursday	Friday
Company $X$	5	4	8	6	2
Company $Y$	3	2	7	4	-1

Table 1: Daily change in share price (£)

Do the following calculations (there is no need to use a calculator):

- Compute the sample covariance between the two sequences to two decimal places.
- Compute the sample correlation between the two sequences. You may leave your answer as a fraction.
- Compute the sample correlation between the two sequences to two decimal places.
- Are the two sequences significantly correlated?

## Solution to Question 2

**Part (a):**

Let  $\mathbf{x} = (5, 4, 8, 6, 2)$  and let  $\mathbf{y} = (3, 2, 7, 4, -1)$ . Then

$$\bar{x} = \frac{1}{5} \sum_{i=1}^5 x_i = \frac{1}{5}(25) = 5$$

$$\bar{y} = \frac{1}{5} \sum_{i=1}^5 y_i = \frac{1}{5}(15) = 3$$

Then,

$$\mathbf{x} - \bar{x} = (0, -1, 3, 1, -3)$$

$$\mathbf{y} - \bar{y} = (0, -1, 4, 1, -4)$$

Which implies that the sample covariance is

$$\begin{aligned} \frac{1}{5-1} \sum_{i=1}^5 (x_i - \bar{x})(y_i - \bar{y}) &= \frac{1}{4} [(0)(0) + (-1)(-1) + (3)(4) + (1)(1) + (-3)(-4)] \\ &= \frac{1}{4} [0 + 1 + 12 + 1 + 12] \\ &= \frac{26}{4} = \frac{13}{2} = 6.5 \end{aligned}$$

**Part (b):**

We need to compute the sample variances of  $\mathbf{x}$  and  $\mathbf{y}$ , or at least the sum of squared deviations:

$$\begin{aligned}\sum_{i=1}^5 (x_i - \bar{x})^2 &= (0)^2 + (1)^2 + (3)^2 + (1)^2 + (-3)^2 = 0 + 1 + 9 + 1 + 9 = 20 \\ \sum_{i=1}^5 (y_i - \bar{y})^2 &= (0)^2 + (-1)^2 + (4)^2 + (1)^2 + (-4)^2 = 0 + 1 + 16 + 1 + 16 = 34\end{aligned}$$

Then the sample correlation is

$$\begin{aligned}r_{XY} &= \frac{\sum_{i=1}^5 (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^5 (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^5 (y_i - \bar{y})^2}} \\ &= \frac{26}{\sqrt{20}\sqrt{34}} = \frac{26}{2\sqrt{5}\sqrt{34}} = \frac{13}{\sqrt{5}\sqrt{34}}\end{aligned}$$

**Part (c):** You can use a calculator for this question, but it is not necessary. Using the computation in Part (b):

$$\begin{aligned}r_{XY} &= \frac{13}{\sqrt{5}\sqrt{34}} = \frac{\sqrt{169}}{\sqrt{170}} = \sqrt{\frac{169}{170}} = \sqrt{\frac{170-1}{170}} = \sqrt{1 - \frac{1}{170}} \\ &> \sqrt{1 - \frac{1}{100}} = \sqrt{0.99} \\ &> 0.99\end{aligned}$$

The first inequality follows since

$$\begin{aligned}170 &> 100 \\ \Rightarrow \frac{1}{170} &< \frac{1}{100} \\ \Rightarrow 1 - \frac{1}{170} &> 1 - \frac{1}{100} \\ \Rightarrow \sqrt{1 - \frac{1}{170}} &> \sqrt{1 - \frac{1}{100}}\end{aligned}$$

The second inequality follows because the functions  $f(x) = \sqrt{x}$  and  $g(x) = x$  on the interval  $(0,1)$  have the property  $f(x) > g(x)$  (to see this, plot the functions).

Finally, although we always have  $r_{XY} \leq 1$ , in this case  $\frac{169}{170} < 1 \Rightarrow \sqrt{\frac{169}{170}} < 1$ , and therefore

$$0.99 < r_{XY} < 1.$$

**Part (d):**

Although we do not have the distribution for  $r_{XY}$  in this case, and so cannot make an inference with any degree of confidence, a value of  $r_{XY} > 0.99$  is likely to be significant.

### Question 3

Suppose that  $X$  and  $Y$  are two normally distributed random variables that are neither independent nor identically distributed. In fact, suppose it is known that  $X \sim N(1, 8)$  and  $Y \sim N(5, 2)$ , and their correlation is  $\text{Cor}(X, Y) = \frac{1}{9}$ . Defining the new random variable  $Z = 2X + Y$ , compute the correlation  $\text{Cor}(Y, Z)$ .

### Solution to Question 3

By properties of the covariance function,

$$\begin{aligned}\text{Cov}(Y, Z) &= \text{Cov}(Y, 2X + Y) \\ &= \text{Cov}(Y, 2X) + \text{Cov}(Y, Y) \\ &= \text{Cov}(2X, Y) + \text{Var}(Y) \\ &= 2\text{Cov}(X, Y) + \text{Var}(Y).\end{aligned}$$

Now, since the correlation is defined as  $\text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}$ , we have

$$\text{Cov}(X, Y) = \sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}\text{Cor}(X, Y) = \sqrt{8}\sqrt{2}\left(\frac{1}{9}\right) = \frac{4}{9}.$$

Then,

$$\begin{aligned}\text{Cov}(Y, Z) &= 2\text{Cov}(X, Y) + \text{Var}(Y) \\ &= 2\left(\frac{4}{9}\right) + 2 \\ &= \frac{8}{9} + 2 \\ &= \frac{26}{9}.\end{aligned}$$

Now, by the properties of the variance (or Exercise 6.1.4 in the notes),

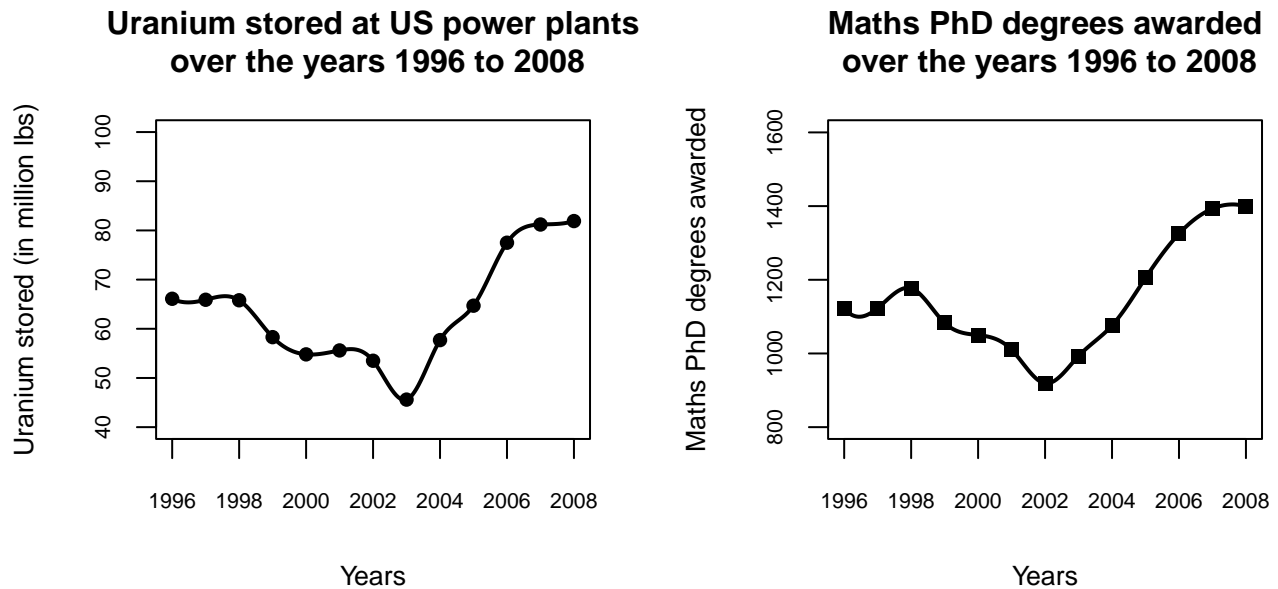
$$\begin{aligned}\text{Var}(Z) &= \text{Var}(2X + Y) \\ &= \text{Var}(2X) + \text{Var}(Y) + 2\text{Cov}(2X, Y) \\ &= 4\text{Var}(X) + \text{Var}(Y) + 4\text{Cov}(X, Y) \\ &= 4 \cdot 8 + 2 + 4\left(\frac{4}{9}\right) \\ &= 34 + \frac{16}{9} \\ &= \frac{322}{9}.\end{aligned}$$

Finally, we can compute the correlation of  $Y$  and  $Z$  as

$$\begin{aligned}\text{Cor}(Y, Z) &= \frac{\text{Cov}(Y, Z)}{\sqrt{\text{Var}(Y)}\sqrt{\text{Var}(Z)}} = \frac{\frac{26}{9}}{\sqrt{2}\sqrt{\frac{322}{9}}} = \frac{\frac{26}{9}}{\frac{2}{3}\sqrt{161}} \\ &= \frac{13}{3\sqrt{161}} \approx 0.3415.\end{aligned}$$

### Question 4

The figures below show  $X$ , the amount of uranium stored in US power plants and  $Y$ , the number of PhD degrees awarded in mathematics in the US for the years 1996 to 2008. The graphs appear to be very similar, and if one computes the sample correlation between the two data sets for  $X$  and  $Y$ , one obtains a value of  $r = 0.952$ . Taking this sample correlation value into account, can we conclude that the two quantities  $X$  and  $Y$  are related and have an influence on each other? Provide justification for your answer.



### Solution to Question 4

We cannot make any conclusion regarding the dependence/influence of  $X$  and  $Y$  on each other, since correlation does not imply causation.