

## Homework 10

Professor: Ziyu Shao &amp; Dingzhu Wen

Due: 23:59 on April 23, 2023

1. Show the proof of general LOTP (four cases).

	$Y$ discrete	$Y$ continuous
$X$ discrete	$P(X = x) = \sum_y P(X = x Y = y)P(Y = y)$	$P(X = x) = \int_{-\infty}^{\infty} P(X = x Y = y)f_Y(y)dy$
$X$ continuous	$f_X(x) = \sum_y f_{X Y}(x y)P(Y = y)$	$f_X(x) = \int_{-\infty}^{\infty} f_{X Y}(x y)f_Y(y)dy$

2. The bus company from Blissville decides to start service in Blotchville, sensing a promising business opportunity. Meanwhile, Fred has moved back to Blotchville. Now when Fred arrives at the bus stop, either of two independent bus lines may come by (both of which take him home). The Blissville company's bus arrival times are exactly 15 minutes apart, whereas the time from one Blotchville company bus to the next is  $\text{Expo}(\frac{1}{15})$ . Fred arrives at a uniformly random time on a certain day.
  - (a) What is the probability that the Blotchville company bus arrives first?
  - (b) What is the CDF of Fred's waiting time for a bus?
3. A chicken lays a  $\text{Pois}(\lambda)$  number  $N$  of eggs. Each egg hatches a chick with probability  $p$ , independently. Let  $X$  be the number which hatch, and  $Y$  be the number which do NOT hatch.
  - (a) Find the joint PMF of  $N, X, Y$ . Are they independent?
  - (b) Find the joint PMF of  $N, X$ . Are they independent?
  - (c) Find the joint PMF of  $X, Y$ . Are they independent?
  - (d) Find the correlation between  $N$  and  $X$ . Your final answer should work out to a simple function of  $p$  and the  $\lambda$  should cancel out.
4. A scientist makes two measurements, considered to be independent standard Normal random variables. Find the correlation between the larger and smaller of the values.

5. This problem explores a visual interpretation of covariance. Data are collected for  $n \geq 2$  individuals, where for each individual two variables are measured (e.g., height and weight). Assume independence across individuals (e.g., person 1's variables gives no information about the other people), but not within individuals (e.g., a person's height and weight may be correlated).

Let  $(x_1, y_1), \dots, (x_n, y_n)$  be the  $n$  data points. The data are considered here as fixed, known numbers—they are the observed values after performing an experiment. Imagine plotting all the points  $(x_i, y_i)$  in the plane, and drawing the rectangle determined by each pair of points. For example, the points  $(1, 3)$  and  $(4, 6)$  determine the rectangle with vertices  $(1, 3), (1, 6), (4, 6), (4, 3)$ .

The signed area contributed by  $(x_i, y_i)$  and  $(x_j, y_j)$  is the area of the rectangle they determine if the slope of the line between them is positive, and is the negative of the area of the rectangle they determine if the slope of the line between them is negative. (Define the signed area to be 0 if  $x_i = x_j$  or  $y_i = y_j$ , since then the rectangle is degenerate.) So the signed area is positive if a higher  $x$  value goes with a higher  $y$  value for the pair of points, and negative otherwise. Assume that the  $x_i$  are all distinct and the  $y_i$  are all distinct.

- (a) The sample covariance of the data is defined to be

$$r = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

where

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \text{ and } \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

are the sample means. (There are differing conventions about whether to divide by  $n - 1$  or  $n$  in the definition of sample covariance, but that need not concern us for this problem.)

Let  $(X, Y)$  be one of the  $(x_i, y_i)$  pairs, chosen uniformly at random. Determine precisely how  $\text{Cov}(X, Y)$  is related to the sample covariance.

- (b) Let  $(X, Y)$  be as in (a), and  $(\tilde{X}, \tilde{Y})$  be an independent draw from the same distribution. That is,  $(X, Y)$  and  $(\tilde{X}, \tilde{Y})$  are randomly chosen from the  $n$  points, independently (so it is possible for the same point to be chosen twice).

Express the total signed area of the rectangles as a constant times  $E((X - \tilde{X})(Y - \tilde{Y}))$ . Then show that the sample covariance of the data is a constant times the total signed area of the rectangles.

Hint: Consider  $E((X - \tilde{X})(Y - \tilde{Y}))$  in two ways: as the average signed area of the random rectangle formed by  $(X, Y)$  and  $(\tilde{X}, \tilde{Y})$ , and using properties of expectation to relate it to  $\text{Cov}(X, Y)$ . For the former, consider the  $n^2$  possibilities for which point  $(X, Y)$  is and which point  $(\tilde{X}, \tilde{Y})$ ; note that  $n$  such choices result in degenerate rectangles.

- (c) Based on the interpretation from (b), give intuitive explanations of why for any r.v.s  $W_1, W_2, W_3$  and constants  $a_1, a_2$ , covariance has the following properties:
- (i)  $\text{Cov}(W_1, W_2) = \text{Cov}(W_2, W_1)$ ;
  - (ii)  $\text{Cov}(a_1 W_1, a_2 W_2) = a_1 a_2 \text{Cov}(W_1, W_2)$ ;
  - (iii)  $\text{Cov}(W_1 + a_1, W_2 + a_2) = \text{Cov}(W_1, W_2)$ ;
  - (iv)  $\text{Cov}(W_1, W_2 + W_3) = \text{Cov}(W_1, W_2) + \text{Cov}(W_1, W_3)$ .