# Introduction to Probability & Statistics

Assignment 4, 2025/26

---

> **ℹ Instructions**
>
> Submit your answers to the **four questions marked ☁ Hand-in**. You should upload your solutions to the VLE as a *single pdf file*. Marks will be awarded for clear, logical explanations, as well as for correctness of solutions.
> Solutions to questions marked ✖ have been released at the same time as this assignment, in case you want to check your answers or need a hint.
> You should also look at the other questions in preparation for your Week 9 seminar.

## Starters

*These questions should help you to gain confidence with the basics.*

**S1.** Let $X$ be a random variable with $\mathbb{E}[X] = 5$. What is the expectation of $3X + 5$? If furthermore $\mathbb{E}[X^2] = 30$, what is the variance of $X$?

> **Answer**
>
> We can use the linearity of expectation to find that $\mathbb{E}[3X + 5] = 3\mathbb{E}[X] + 5 = 20$. The variance is $\text{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = 30 - 5^2 = 5$.

**S2.** ✖ I arrive at the train station at 12.00 exactly. My train departs at a time which follows a (continuous) uniform distribution on the interval [11.55, 12.15]. What is the probability that I miss my train?

> **Answer**
>
> Let $X$ denote the random time after 11.55 at which the train leaves. The question tells us that $X \sim \text{Uniform}[0, 20]$. I miss the train if $X < 5$, which has probability
>
> $$\mathbb{P}(X < 5) = \int_0^5 \frac{1}{20} dx = \frac{1}{4}.$$

**S3.** ☁ **Hand-in**

Buses leave campus for the train station every 20 minutes, at 5, 25, and 45 minutes past the hour. If a student arrives at the bus stop at a time that follows a (continuous) uniform distribution on the interval between 09.00 and 09.35, find the probability that they wait

    a) less than 5 minutes for a bus;
    b) at least 10 minutes for a bus.

**S4.** Suppose that you have a lecture at 14.00, and that the time taken to travel from your room to the lecture theatre is normally distributed with mean 30 minutes and standard deviation 4 minutes. What is the latest time you should leave your room if you want to be 99% certain that you will not miss the start of the lecture? (Hint: if $Z \sim \mathsf{N}(0,1)$ then the R function `qnorm(p)` returns the value $z \in \mathbb{R}$ such that $\mathbb{P}\left(Z \leq z\right) = p$.)

**S5.** ⚒ Give an example of a joint probability table for two discrete random variables $X$ and $Y$, each having only two possible values, so that $F_{X,Y}(5, 6) = 0.4, F_X(5) = 0.5, F_Y(6) = 0.6$ and $\mathbb{E}\left[X\right] = 10, \mathbb{E}\left[Y\right] = 4$.

**S6.** ⚒ Let $X : \Omega \to \{1, 2\}$ and $Y : \Omega \to \{0, 1\}$ be two discrete random variables. The following is a partial table of their joint and their marginal mass functions:

| $y \backslash x$ | 1 | 2 | $p_Y(y)$ |
|---|---|---|---|
| 0 | 1/6 | 1/2 | |

| $y \backslash x$ | 1 | 2 | $p_Y(y)$ |
|---|---|---|---|
| 1 | | | |
| $p_X(x)$ | 5/12 | | 1 |

a) Fill in the missing values.
b) Determine the joint distribution function of $X$ and $Y$.
c) Calculate $\mathbb{E}[X]$ and $\mathbb{E}[Y]$.
d) Let $Z = XY$. Calculate $\mathbb{E}[Z]$.

---

**Answer**

a) The missing entries in the probability table are determined by the requirement that summing the joint probabilities across a row or across a column in the table gives the corresponding marginal probability and by the requirement that the marginal probabilities for $X$ as well as those for $Y$ have to add up to 1. So first we determine $p_Y(0) = 1/6 + 1/2 = 2/3$. Then we can determine $p_Y(1) = 1 - p_Y(0) = 1 - 2/3 = 1/3$ and $p_X(2) = 1 - p_X(1) = 1 - 5/12 = 7/12$. Finally we determine $p_{X,Y}(1,1) = p_X(1) - p_{X,Y}(1,0) = 5/12 - 1/6 = 1/4$ and $p_{X,Y}(2,1) = p_X(2) - p_{X,Y}(2,0) = 7/12 - 1/2 = 1/12$.

| $y \backslash x$ | 1 | 2 | $p_Y(y)$ |
|---|---|---|---|
| 0 | 1/6 | 1/2 | 2/3 |
| 1 | 1/4 | 1/12 | 1/3 |
| $p_X(x)$ | 5/12 | 7/12 | 1 |

b) The joint distribution function $F_{X,Y}(x, y)$ is by definition given by $\mathbb{P}(X \leq x, Y \leq y)$. So for example

$$F_{X,Y}(1.5, 1.5) = p_{X,Y}(1, 0) + p_{X,Y}(1, 1) = \frac{1}{6} + \frac{1}{4} = \frac{5}{12}.$$

By doing more such calculations we find that

$$F_{X,Y} = \begin{cases} 0 & \text{if } x < 1 \text{ or } y < 0 \\ 1/6 & \text{if } x \in [1, 2) \text{ and } y \in [0, 1) \\ 5/12 & \text{if } x \in [1, 2) \text{ and } y \geq 1 \\ 2/3 & \text{if } x \geq 2 \text{ and } y \in [0, 1) \\ 1 & \text{if } x \geq 2 \text{ and } y \geq 1. \end{cases}$$

c) For calculating the expectations of $X$ and $Y$ we can use their marginal mass functions:

$$\mathbb{E}[X] = 1 \cdot p_X(1) + 2 \cdot p_X(2) = 1 \cdot \frac{5}{12} + 2 \cdot \frac{7}{12} = \frac{19}{12}$$

and

$$\mathbb{E}[Y] = 0 \cdot p_Y(0) + 1 \cdot p_Y(1) = p_Y(1) = \frac{1}{3}.$$

d) The random variable $Z = XY$ can take the possible values 0, 1 and 2 with probabilities

$$p_Z(0) = p_{X,Y}(1, 0) + p_{X,Y}(2, 0) = p_Y(0) = \frac{2}{3}$$

$$p_Z(1) = p_{X,Y}(1, 1) = \frac{1}{4}, \quad p_Z(2) = p_{X,Y}(2, 1) = \frac{1}{12}.$$

Thus

$$\mathbb{E}[Z] = 1 \cdot p_Z(1) + 2 \cdot p_Z(2) = \frac{1}{4} + 2\frac{1}{12} = \frac{5}{12}.$$

## S7. ☁ Hand-in

Let $X : \Omega \to \{0, 1\}$ and $Y : \Omega \to \{0, 1\}$ be two discrete random variables. The following is a partial table of their joint and their marginal mass functions:

| $y \backslash x$ | 0 | 1 | $p_Y(y)$ |
|---|---|---|---|
| 0 | 1/8 | | 1/4 |
| 1 | | 1/2 | |
| $p_X(x)$ | | | |

a) Fill in the missing values.
b) Calculate $\mathbb{E}[XY]$.

---

**Answer**

a) The completed table is as follows **[2 marks]**

| $y \backslash x$ | 0 | 1 | $p_Y(y)$ |
|---|---|---|---|
| 0 | 1/8 | 1/8 | 1/4 |
| 1 | 1/4 | 1/2 | 3/4 |
| $p_X(x)$ | 3/8 | 5/8 | 1 |

b) Note that the random variable $XY$ equals 0 unless $(X, Y) = (1, 1)$, in which case $XY = 1$. (This means that $XY$ has a Bernoulli distribution.) Thus $\mathbb{E}[XY] = \mathbb{P}((X, Y) = (1, 1)) = 1/2$. **[3 marks]**

---

# Mains

*These are important, and cover some of the most substantial parts of the course.*

**M1.** ✖ A random variable $W$ has probability density function

$$f_W(x) = \begin{cases} \frac{6}{5675}(5x^2 + 3x + 11) & \text{for } 3 \leq x \leq 8 \\ 0 & \text{otherwise.} \end{cases}$$

Would you expect $\mathbb{E}[W]$ to lie closer to 3 or to 8? Calculate $\mathbb{E}[W]$ and check whether your intuition was correct.

> **Answer**
>
> Since $f_W$ is increasing on the interval $[3, 8]$ we know from the interpretation of expectation as centre of mass that the expectation should lie closer to 8 than to 3. The computation:
>
> $$\mathbb{E}[W] = \int_3^8 x f_W(x) dx = \frac{6}{5675} \int_3^8 \left(5x^3 + 3x^2 + 11x\right) dx = \frac{2787}{454} = 6.14.$$

**M2.** Let $X \sim \text{Geom}(p)$. Calculate $\mathbb{E}[h(X)]$, where $h(x) = e^{tx}$ for some $t > 0$. For what values of $t$ is $\mathbb{E}[h(X)] < \infty$?

> **Answer**
>
> We use the formula for the expectation of a function of a discrete random variable:
>
> $$\mathbb{E}[h(X)] = \sum_{k=1}^{\infty} h(k)p(1-p)^{k-1} = \sum_{k=1}^{\infty} e^{tk}p(1-p)^{k-1}$$
> $$= pe^t \sum_{k=1}^{\infty} \left[e^t(1-p)\right]^{k-1} = pe^t \sum_{k=0}^{\infty} \left[e^t(1-p)\right]^{k}$$
> $$= \frac{pe^t}{1 - e^t(1-p)}.$$
>
> This final step requires $e^t(1-p) < 1$. (Otherwise the geometric sum does not converge to a finite limit.)

**M3.** ☁ **Hand-in**

Show that if $Z$ is a standard normal random variable then, for $x > 0$,

a) $\mathbb{P}(Z > x) = \mathbb{P}(Z < -x)$;
b) $\mathbb{P}(|Z| > x) = 2\mathbb{P}(Z > x)$;
c) $\mathbb{P}(|Z| < x) = 2\mathbb{P}(Z < x) - 1$.

Hint: for part (a), express the probabilities in terms of integrals over the density function $\phi$, and use the fact that $\phi$ is an even function (i.e. $\phi(z) = \phi(-z)$).

There are many ways to show these identities. We use the hint about the symmetry of the density function of a standard normal random variable:

$$\phi(-z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(-z)^2}{2}\right) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) = \phi(z).$$

a)

$$\mathbb{P}\left(Z > x\right) = \int_x^\infty \phi(z)dz = \int_{-\infty}^{-x} \phi(-u)du = \int_{-\infty}^{-x} \phi(u)du = \mathbb{P}\left(Z < -x\right);$$

**[1 mark]**

b)

$$\mathbb{P}\left(|Z| > x\right) = \mathbb{P}\left(Z > x\right) + \mathbb{P}\left(Z < -x\right) = 2\mathbb{P}\left(Z > x\right),$$

where the last equality follows from part (a). **[2 marks]**

c)

$$\mathbb{P}\left(|Z| < x\right) = 1 - \mathbb{P}\left(|Z| > x\right) = 1 - 2\mathbb{P}\left(Z > x\right)$$
$$= 1 - 2\left(1 - \mathbb{P}\left(Z < x\right)\right) = 2\mathbb{P}\left(Z < x\right) - 1,$$

where the second equality follows from part (b). **[2 marks]**

**M4.** ☁ **Hand-in**

Let $X \sim \text{Exp}(\lambda)$. Use proof by induction to show that

$$\mathbb{E}\left[X^m\right] = \frac{m!}{\lambda^m} \qquad \text{for all } m \in \mathbb{N} \cup \{0\}.$$

The statement is true for $m = 0$:

$$\mathbb{E}\left[X^0\right] = \mathbb{E}\left[1\right] = 1 = \frac{0!}{\lambda^0}.$$

**[1 mark]**

Now suppose that the statement holds for some $k \in \mathbb{N} \cup \{0\}$, and consider the case $m = k + 1$:

$$\begin{aligned}
\mathbb{E}\left[X^{k+1}\right] &= \int_{-\infty}^\infty x^{k+1} f_X(x)dx = \int_0^\infty x^{k+1} \lambda e^{-\lambda x}dx \\
&= \int_0^\infty x^{k+1} \frac{d}{dx}\left(-e^{-\lambda x}\right)dx \\
&= -\left[x^{k+1}e^{-\lambda x}\right]_0^\infty + \int_0^\infty (k+1)x^k e^{-\lambda x}dx \\
&= 0 + \frac{k+1}{\lambda} \int_0^\infty x^k \lambda e^{-\lambda x}dx \\
&= \frac{k+1}{\lambda} \mathbb{E}\left[X^k\right] = \frac{k+1}{\lambda} \frac{k!}{\lambda^k} \quad \text{by our induction hypothesis} \\
&= \frac{(k+1)!}{\lambda^{k+1}}.
\end{aligned}$$

Thus the statement holds for all $m \in \mathbb{N} \cup \{0\}$ by induction. **[4 marks]**

**M5.** Let $X$ and $Y$ be random variables. Show that $\text{Cov}\left(X, Y\right) = \mathbb{E}\left[XY\right] - \mathbb{E}\left[X\right]\mathbb{E}\left[Y\right]$.

We start from the definition of covariance, and use linearity of expectation:

$$\begin{aligned}
\text{Cov}\,(X,Y) &= \mathbb{E}\left[(X - \mathbb{E}\,[X])\,(Y - \mathbb{E}\,[Y])\right] \\
&= \mathbb{E}\,[XY - X\,\mathbb{E}\,[Y] - \mathbb{E}\,[X]\,Y + \mathbb{E}\,[X]\,\mathbb{E}\,[Y]] \\
&= \mathbb{E}\,[XY] - \mathbb{E}\,[X]\,\mathbb{E}\,[Y] - \mathbb{E}\,[X]\,\mathbb{E}\,[Y] + \mathbb{E}\,[X]\,\mathbb{E}\,[Y] \\
&= \mathbb{E}\,[XY] - \mathbb{E}\,[X]\,\mathbb{E}\,[Y].
\end{aligned}$$

**M6.** ⚒ The joint probability mass function $p_{X,Y}(x,y)$ of two random variables $X$ and $Y$ is summarised by the following table, where $\eta$ is some real number:

| $x\backslash y$ | -1 | 0 | 1 |
|---|---|---|---|
| 4 | $\eta - 1/16$ | $1/4 - \eta$ | 0 |
| 5 | $1/8$ | $3/16$ | $1/8$ |
| 6 | $\eta + 1/16$ | $1/16$ | $1/4 - \eta$ |

a) Extend the table by including also the marginal probabilities, i.e., the values of the probability mass functions $p_X$ and $p_Y$.

b) Which are the valid choices for $\eta$?

c) Is there a value of $\eta$ for which $X$ and $Y$ are independent?

a) We extend the probability table to also include the marginal probability mass functions $p_X$ and $p_Y$:

| $x\backslash y$ | -1 | 0 | 1 | $p_X(x)$ |
|---|---|---|---|---|
| 4 | $\eta - 1/16$ | $1/4 - \eta$ | 0 | $3/16$ |
| 5 | $1/8$ | $3/16$ | $1/8$ | $7/16$ |
| 6 | $\eta + 1/16$ | $1/16$ | $1/4 - \eta$ | $3/8$ |
| $p_Y(y)$ | $2\eta + 1/8$ | $1/2 - \eta$ | $3/8 - \eta$ | 1 |

b) All entries of the probability table must be non-negative and they must sum up to $1$. In order for $p_{X,Y}(4,-1)$ to be non-negative we need $\eta \geq 1/16$. In order for $p_{X,Y}(4,0)$ and $p_{X,Y}(6,1)$ to be non-negative we need $\eta \leq 1/4$. The sum over all entries is not affected by the value of $\eta$, so does not give any additional constraints. Therefore any $\eta \in [1/16, 1/4]$ is a valid choice.

c) It is easy to find counterexamples to the factorisation of the joint probability mass function that would have to hold if $X$ and $Y$ were independent. For example

$$p_X(4)p_Y(1) = \frac{3}{16}\left(\frac{3}{8} - \eta\right) \neq 0 = p_{X,Y}(4,1)$$

unless $\eta = 3/8$. However the value $\eta = 3/8$ is not allowed, and hence $X$ and $Y$ can never be independent.

**M7.** A married couple decide to have children until they have at least one child of each sex: let $X$ denote the total number of children that they have. The probability of any one child being a boy is $1/2$

(with the sex of each child being independent of all the others).

a) What is the mass function of $X$? (I.e. write down $\mathbb{P}(X = n)$ for all $n \in \mathbb{N}$.)

b) Show that
$$\mathbb{E}[X] = 3.$$

Hint: you may find it useful to refer to the result from lectures that if $Y \sim \text{Geom}(p)$ then $\mathbb{E}[Y] = 1/p$.

---

**Answer**

a) Clearly the couple need to have at least two children, so $\mathbb{P}(X = 1) = 0$. For $n \geq 2$, there are two ways in which the couple can have exactly $n$ children: either they have $n - 1$ boys in a row, and then a girl; or they have $n - 1$ girls and then a boy. Each of these possibilities has probability $(1/2)^n$. Thus
$$\mathbb{P}(X = n) = (1/2)^n + (1/2)^n = (1/2)^{n-1}, \qquad n \geq 2.$$

b) Here are two possible ways of calculating $\mathbb{E}[X]$.

**Method 1:** We use the usual formula for the expectation of a discrete random variable:
$$\mathbb{E}[X] = \sum_{n=2}^{\infty} n\mathbb{P}(X = n) = \sum_{n=2}^{\infty} n(1/2)^{n-1}$$

Using the hint, we know that if $Y \sim \text{Geom}(1/2)$ then $\mathbb{E}[Y] = 2$. That is,
$$\sum_{n=1}^{\infty} n(1/2)(1/2)^{n-1} = 2.$$

We now manipulate our expression for the expectation, until it looks like something involving this result:
$$\mathbb{E}[X] = \sum_{n=2}^{\infty} n(1/2)^{n-1} = 2\sum_{n=2}^{\infty} n(1/2)(1/2)^{n-1}$$
$$= 2\left[\sum_{n=1}^{\infty} n(1/2)(1/2)^{n-1} - 1/2\right]$$
$$= 2[2 - 1/2] = 3.$$

**Method 2:** If we let $Y = X - 1$ then this random variable takes values in the set $\mathbb{N}$ and has mass function
$$\mathbb{P}(Y = n) = \mathbb{P}(X - 1 = n) = \mathbb{P}(X = n + 1) = (1/2)^n$$
for $n \in \mathbb{N}$. Thus $Y \sim \text{Geom}(1/2)$. It follows that $\mathbb{E}[X] = \mathbb{E}[Y + 1] = \mathbb{E}[Y] + 1 = 2 + 1 = 3$. (Here we're effectively observing that the couple start by having one child, who could be of either sex; they then need to have an additional random number of children until they have one of the opposite sex to the first – this is like repeating independent Bernoulli trials, with "success" meaning that they have a child of the opposite sex.)

---

**M8.** ⚒ Let $X$ be a discrete random variable. Show that for all functions $h_1, h_2 : \mathbb{R} \to \mathbb{R}$,
$$\mathbb{E}[h_1(X) + h_2(X)] = \mathbb{E}[h_1(X)] + \mathbb{E}[h_2(X)].$$

**M9.** Let $X$ and $Y$ be random variables and let $r, s, t, u \in \mathbb{R}$. Show that

$$\rho(rX + s, tY + u) = \begin{cases} \rho(X, Y) & \text{if } rt > 0 \\ 0 & \text{if } rt = 0 \\ -\rho(X, Y) & \text{if } rt < 0 \end{cases}$$

where $\rho(X, Y)$ denotes the correlation coefficient of $X$ and $Y$.

> **Answer**
>
> Let us first assume that $\text{Var}\left(X\right)\text{Var}\left(Y\right) > 0$ and $rt > 0$. Then the definition of the correlation coefficient gives
>
> $$\rho(rX + s, tY + u) = \frac{\text{Cov}\left(rX + s, tY + u\right)}{\sqrt{\text{Var}\left(rX + s\right)\text{Var}\left(tY + u\right)}}. \tag{1}$$
>
> We already know that
>
> $$\text{Var}\left(rX + s\right) = r^2\text{Var}\left(X\right), \quad \text{Var}\left(tY + u\right) = t^2\text{Var}\left(Y\right). \tag{2}$$
>
> We need to derive a similar transformation rule for the covariance.
>
> $$\begin{aligned} \text{Cov}\left(rX + s, tY + u\right) &= \mathbb{E}\left[(rX + s - \mathbb{E}\left[rX + s\right])(tY + u - \mathbb{E}\left[tY + u\right])\right] \\ &= \mathbb{E}\left[(rX + s - (r\mathbb{E}\left[X\right] + s))(tY + u - (t\mathbb{E}\left[Y\right] + u))\right] \\ &= \mathbb{E}\left[r\left(X - \mathbb{E}\left[X\right]\right)t\left(Y - \mathbb{E}\left[Y\right]\right)\right] \\ &= rt\mathbb{E}\left[\left(X - \mathbb{E}\left[X\right]\right)\left(Y - \mathbb{E}\left[Y\right]\right)\right] \\ &= rt\text{Cov}\left(X, Y\right) , \end{aligned} \tag{3}$$
>
> where we repeatedly used the linearity of expectation. Using the transformation rules Equation 2 and Equation 3 in Equation 1 gives
>
> $$\rho(rX + s, tY + u) = \frac{rt}{\sqrt{r^2t^2}}\frac{\text{Cov}\left(X, Y\right)}{\sqrt{\text{Var}\left(X\right)\text{Var}\left(Y\right)}}.$$
>
> The statement now follows from the observation that
>
> $$\frac{rt}{\sqrt{r^2t^2}} = \begin{cases} 1 & \text{if } rt > 0 \\ -1 & \text{if } rt < 0. \end{cases}$$
>
> In case $\text{Var}\left(X\right)\text{Var}\left(Y\right) = 0$ or $rt = 0$ also $\text{Var}\left(rX + s\right)\text{Var}\left(tY + u\right) = rt\text{Var}\left(X\right)\text{Var}\left(Y\right) = 0$, and thus $\rho(rX + s, tY + u) = 0$ by definition. This agrees with the statement because when $\text{Var}\left(X\right)\text{Var}\left(Y\right) = 0$ also $\rho(X, Y) = 0$.

# Desserts

*Still hungry for more? Try these if you want to push yourself further. (These are mostly harder than I'd expect you to answer in an exam, or involve non-examinable material.)*

**D1.** Prove that binomial coefficients satisfy the identity

$$n\binom{n-1}{r-1} = r\binom{n}{r}.$$

Use this to find $\mathbb{E}[X]$ and $\text{Var}(X)$, where $X \sim \text{Bin}(n,p)$.

---

**Answer**

First we prove the identity:

$$n\binom{n-1}{r-1} = n\frac{(n-1)!}{(r-1)!(n-r)!} = r\frac{n!}{r!(n-r)!} = r\binom{n}{r}.$$

For the mean and variance, remember that, since $p_X(\cdot)$ is a mass function, it must sum to one. That is,

$$\sum_{k=0}^{n}\binom{n}{k}p^k(1-p)^{n-k} = 1. \tag{4}$$

Now,

$$\begin{aligned}
\mathbb{E}[X] &= \sum_{k=0}^{n} k\binom{n}{k}p^k(1-p)^{n-k}\\
&= \sum_{k=1}^{n} n\binom{n-1}{k-1}p^k(1-p)^{n-k} \quad \text{(by our identity)}\\
&= np\sum_{k=1}^{n}\binom{n-1}{k-1}p^{k-1}(1-p)^{n-k}\\
&= np\sum_{j=0}^{n-1}\binom{n-1}{j}p^j(1-p)^{(n-1)-j} \quad \text{(putting } j=k-1)\\
&= np,
\end{aligned}$$

thanks to Equation 4.

Furthermore,

$$\begin{aligned}
\mathbb{E}[X(X-1)] &= \sum_{k=0}^{n} k(k-1)\binom{n}{k}p^k(1-p)^{n-k}\\
&= \sum_{k=0}^{n} k(k-1)\frac{n!}{k!(n-k)!}p^k(1-p)^{n-k}\\
&= \sum_{k=0}^{n} \frac{n!}{(n-k)!(k-2)!}p^k(1-p)^{n-k}\\
&= n(n-1)p^2\sum_{k=2}^{n}\frac{(n-2)!}{((n-2)-(k-2))!(k-2)!}p^{k-2}(1-p)^{(n-2)-(k-2)}\\
&= n(n-1)p^2\sum_{j=0}^{n-2}\binom{n-2}{j}p^j(1-p)^{(n-2)-j} \quad \text{(putting } j=k-2)\\
&= n(n-1)p^2,
\end{aligned}$$

again thanks to Equation 4. It follows that

$$\mathbb{E}\left[X^2\right] = n(n-1)p^2 + np\,,$$

and so

$$\text{Var}\left(X\right) = \mathbb{E}\left[X^2\right] - \mathbb{E}\left[X\right]^2 = np(1-p)\,.$$

**D2.** Let $X \sim \text{Uniform}(0, a)$ for some $a > 0$. Show that for any $n \in \mathbb{N}$,

$$\mathbb{E}\left[X^n\right] = \frac{a^n}{n+1}\,.$$

Use this to determine $\rho(X, X^2)$, and show that this does not depend upon the value of $a$.

> **Answer**
>
> For $n \in \mathbb{N}$ we calculate
>
> $$\mathbb{E}\left[X^n\right] = \int_{-\infty}^{\infty} x^n f_X(x)dx = \int_0^a \frac{x^n}{a}dx = \frac{1}{a}\left[\frac{x^{n+1}}{n+1}\right]_0^a = \frac{a^n}{n+1}\,.$$
>
> Now we calculate the covariance of $X$ and $X^2$:
>
> $$\text{Cov}\left(X, X^2\right) = \mathbb{E}\left[X^3\right] - \mathbb{E}\left[X\right]\mathbb{E}\left[X^2\right] = \frac{a^3}{4} - \frac{a^2}{3}\frac{a}{2} = \frac{a^3}{12}\,.$$
>
> We also have
>
> $$\text{Var}\left(X\right) = \mathbb{E}\left[X^2\right] - \mathbb{E}\left[X\right]^2 = \frac{a^2}{3} - \left(\frac{a}{2}\right)^2 = \frac{a^2}{12}$$
>
> and
>
> $$\text{Var}\left(X^2\right) = \mathbb{E}\left[X^4\right] - \mathbb{E}\left[X^2\right]^2 = \frac{a^4}{5} - \left(\frac{a^2}{3}\right)^2 = \frac{4a^4}{45}\,.$$
>
> Finally, we calculate
>
> $$\rho(X, X^2) = \frac{\text{Cov}\left(X, X^2\right)}{\sqrt{\text{Var}\left(X\right)\text{Var}\left(X^2\right)}} = \frac{a^3/12}{\sqrt{a^6/135}} = \frac{\sqrt{135}}{12} = \frac{\sqrt{15}}{4}\,,$$
>
> which doesn't depend upon $a$.

**D3.** A bag contains 3 cubes, 4 pyramids and 7 spheres. An object is drawn randomly from the bag and its type is recorded. Then the object is replaced. This is repeated 20 times.

  a. Let $C_i$ be the indicator random variable for the event that the $i$-th draw gives a cube, for $i = 1, \ldots, 20$. Calculate $\mathbb{E}\left[C_i\right]$, $\mathbb{E}\left[C_i^2\right]$ and $\mathbb{E}\left[C_i C_j\right]$ for $i \neq j$.

  b. Let $C$ be the number of times a cube was drawn, Use that $C = \sum_{i=1}^{20} C_i$ to calculate $\mathbb{E}\left[C\right]$ and $\text{Var}\left(C\right)$.

  c. Let $S_i$ be the indicator random variable for the event that the $i$-th draw gives a sphere. Calculate $\mathbb{E}\left[C_i S_i\right]$ and $\mathbb{E}\left[C_i S_j\right]$ for $i \neq j$.

  d. Let $S$ be the number of times a sphere was drawn. Use the above results to calculate $\mathbb{E}\left[CS\right]$, $\text{Cov}\left(C, S\right)$, $\rho(C, S)$.

a. As three of the 14 shapes are cubes, the probability to draw a cube is $3/14$. Hence $C_i \sim \mathsf{Bern}(3/14)$. This immediately gives

$$\mathbb{E}\left[C_i\right] = \mathbb{E}\left[C_i^2\right] = \frac{3}{14}.$$

For $i \neq j$ the event that the $i$-th draw gives a cube and the event that the $j$-th cube gives a draw are independent (because we put the shape back after each draw). Thus the indicator random variables $C_i$ and $C_j$ for these events are also independent and thus

$$\mathbb{E}\left[C_i C_j\right] = \mathbb{E}\left[C_i\right] \mathbb{E}\left[C_j\right] = \left(\frac{3}{14}\right)^2 = \frac{9}{196}.$$

b. The linearity of expectation gives

$$\mathbb{E}\left[C\right] = \mathbb{E}\left[\sum_{i=1}^{20} C_i\right] = \sum_{i=1}^{20} \mathbb{E}\left[C_i\right] = 20\frac{3}{14} = \frac{30}{7}.$$

Because the $C_i$ are independent of each other, the variance of their sum equals the sum of their variances:

$$\mathsf{Var}\left(C\right) = \mathsf{Var}\left(\sum_{i=1}^{20} C_i\right) = \sum_{i=1}^{20} \mathsf{Var}\left(C_i\right) = 20\frac{3}{14}\frac{11}{14} = \frac{165}{49}.$$

c. We observe that $C_i S_i = 0$ because on the same draw one can not simultaneously have a cube and a sphere. Thus also $\mathbb{E}\left[C_i S_i\right] = 0$. If $i \neq j$ we can use independence to factorise the expectation:

$$\mathbb{E}\left[C_i S_j\right] = \mathbb{E}\left[C_i\right] \mathbb{E}\left[S_j\right] = \frac{3}{14}\frac{1}{2} = \frac{3}{28},$$

where we used that the probability of drawing a sphere is $1/2$.

d. We have

$$\mathbb{E}\left[CS\right] = \mathbb{E}\left[\sum_{i=1}^{20} C_i \sum_{j=1}^{20} S_j\right] = \sum_{i=1}^{20}\sum_{j=1}^{20} \mathbb{E}\left[C_i S_j\right].$$

We split the sum over all pairs $(i,j)$ into the pairs where $i \neq j$ and the pairs $(i,i)$, so

$$\mathbb{E}\left[CS\right] = \sum_{i=1}^{20}\sum_{\substack{j=1 \\ j\neq i}}^{20} \mathbb{E}\left[C_i S_j\right] + \sum_{i=1}^{20} \mathbb{E}\left[C_i S_i\right].$$

Using our above results for $\mathbb{E}\left[C_i S_i\right]$ and $\mathbb{E}\left[C_i S_j\right]$ and recognising that there are $20 \cdot 19 = 380$ pairs where $i \neq j$ this gives us

$$\mathbb{E}\left[CS\right] = \sum_{i=1}^{20}\sum_{\substack{j=1 \\ j\neq i}}^{20} \frac{3}{28} + \sum_{i=1}^{20} 0 = 380\frac{3}{28} = \frac{285}{7}.$$

We also calculate

$$\mathbb{E}\left[S\right] = \mathbb{E}\left[\sum_{i=1}^{20} S_i\right] = \sum_{i=1}^{20} \mathbb{E}\left[S_i\right] = 20\frac{1}{2} = 10.$$

The covariance can then be calculated as

$$\mathsf{Cov}\left(C, S\right) = \mathbb{E}\left[CS\right] - \mathbb{E}\left[C\right]\mathbb{E}\left[S\right] = \frac{285}{7} - \frac{30}{7}10 = -\frac{15}{7}.$$

To calculate the correlation coefficient we also need

$$\text{Var}\,(S) = \text{Var}\left(\sum_{i=1}^{20} S_i\right) = \sum_{i=1}^{20} \text{Var}\,(S_i) = 20\frac{1}{2}\frac{1}{2} = 5.$$

The correlation coefficient is

$$\rho(C,S) = \frac{\text{Cov}\,(C,S)}{\sqrt{\text{Var}\,(C)\,\text{Var}\,(S)}} = -\sqrt{\frac{3}{11}} \approx -0.5222.$$

**D4.** Consider a random variable $X \sim \text{Uniform}[a,b]$, where $a$ and $b$ are unknown. You are told that

$$\mathbb{P}\,(X < 2) = 1/3 \quad \text{and} \quad \mathbb{P}\,(1 < X \le 3) = 1/2\,.$$

Given this information, find $a$ and $b$.

---

**Answer**

From the first equation we immediately know that $a < 2 < b$. Now, for a continuous random variable, we obtain the probability that it lies in an interval $(c,d)$ by integrating the density function over that interval, i.e.

$$\mathbb{P}\,(c \le X \le d) = \int_c^d f_X(x)dx\,.$$

Since $X \sim \text{Uniform}[a,b]$, we know that

$$f_X(x) = \begin{cases} 1/(b-a) & x \in [a,b] \\ 0 & \text{otherwise.} \end{cases}$$

Thus we obtain

$$1/3 = \mathbb{P}\,(X < 2) = \mathbb{P}\,(a \le X < 2) = \int_a^2 1/(b-a)dx = (2-a)/(b-a)\,. \tag{5}$$

In order to use the second equation ($\mathbb{P}\,(1 < X \le 3) = 1/2$) in the same way, we have two possibilities to consider:

1. $a < 1$
2. $1 \le a < 2$

Suppose first that $a < 1$. Then

$$1/2 = \mathbb{P}\,(1 < X \le 3) = \int_1^3 1/(b-a)dx = 2/(b-a)\,, \tag{6}$$

since the density function $f_X$ is equal to $1/(b-a)$ for all $x \in [1,3]$ if $a < 1$.
If $1 \le a$ however, then instead we obtain

$$1/2 = \mathbb{P}\,(1 < X \le 3) = \int_1^a 0\,dx + \int_a^3 1/(b-a)dx = (3-a)/(b-a)\,. \tag{7}$$

We now have to solve these simultaneous equations in order to find $a$ and $b$. If we assume that $1 \le a < 2$, then we must try to solve Equation 5 and Equation 7 together; but this gives

$$2(3-a) = 3(2-a)\,,$$

resulting in $a = 0$. But this contradicts our assumption that $1 \le a$!
So it must be the case that $a < 1$: now we must solve Equation 5 and Equation 6, and this *is* possible, with $a = 2/3$ and $b = 14/3$.

**D5.** Let $X$ and $Y$ be two independent geometrically distributed random variables with parameter $p$, i.e., $X \sim \text{Geom}(p)$ and $Y \sim \text{Geom}(p)$. For any natural numbers $i$ and $n$ with $i < n$ calculate the conditional probability $\mathbb{P}\left(X = i \,|\, X + Y = n\right)$. Describe in words the meaning in terms of Bernoulli trials of what you just calculated.

---

**Answer**

According to the definition of conditional probability,

$$\mathbb{P}\left(X = i \,|\, X + Y = n\right) = \frac{\mathbb{P}\left(X = i,\, X + Y = n\right)}{\mathbb{P}\left(X + Y = n\right)}.$$

For the numerator we can use that the event $\{X = i, X + Y = n\}$ is the event $\{X = i, Y = n - i\}$. We then know that the independence of $X$ and $Y$ implies the factorisation of that probability:

$$\mathbb{P}\left(X = i,\, X + Y = n\right) = \mathbb{P}\left(X = i, Y = n - i\right) = \mathbb{P}\left(X = i\right)\mathbb{P}\left(Y = n - i\right).$$

We can now substitute in the probability mass function for the geometric distribution with parameter $p$:

$$\mathbb{P}\left(X = i\right) = (1 - p)^{i-1} p$$

and thus

$$\mathbb{P}\left(Y = n - i\right) = (1 - p)^{n-i-1} p.$$

This gives

$$\mathbb{P}\left(X = i,\, X + Y = n\right) = (1 - p)^{i-1} p (1 - p)^{n-i-1} p = (1 - p)^{n-2} p^2.$$

Note that this is independent of $i$.
For the denominator we use the partition theorem to write

$$\mathbb{P}\left(X + Y = n\right) = \sum_{i=1}^{n-1} \mathbb{P}\left(X = i, X + Y = n\right).$$

From our calculation above we see that every term in the sum is the same, so

$$\mathbb{P}\left(X + Y = n\right) = (n - 1)\mathbb{P}\left(X = i, X + Y = n\right).$$

Putting this all together we finally find that

$$\mathbb{P}\left(X = i \,|\, X + Y = n\right) = \frac{\mathbb{P}\left(X = i, X + Y = n\right)}{\mathbb{P}\left(X + Y = n\right)} = \frac{1}{n - 1}.$$

A geometric random variable counts the number of turns until the first success in repeated Bernoulli trials. Therefore the sum $X + Y$ of two identical and independent geometric random variables counts the number of turns until the *second* success. So the conditional probability we calculated is the probability that the first success happens on a particular trial $i$ given that the second success happens on the $n$-th trial. The result shows that the first success is then equally likely to occur on any of the $n - 1$ trials before the $n$-th trial.