

Probability distributions EBP038A05

Study Guide

Nicky D. van Foreest, Ruben van Beesten
Joost Doornbos, Wietze Koops, Mikael Makonnen, Zexuan Yuan

January 20, 2022

CONTENTS

| | | |
|------|--|----|
| 7 | Questions and remarks | 7 |
| 7.1 | Simple questions | 7 |
| 7.2 | Memoryless excursions: A confusing problem with memoryless rvs | 10 |
| 7.3 | Exercises on 2D integration | 17 |
| 7.4 | BH exercises: hints and solutions | 19 |
| 7.5 | Challenge: A uniqueness property of the Poisson distribution | 20 |
| 7.6 | Challenge: Improper integrals and the Cauchy distribution | 20 |
| 7.7 | Challenge: Proof about independence of normal rvs | 20 |
| 7.8 | Challenge: Recourse models | 21 |
| 8 | Chapter 8: Questions and remarks | 23 |
| 8.1 | Simple questions | 23 |
| 8.2 | BH exercises: hints and solutions | 28 |
| 8.3 | Challenge: Ping pong balls in a Beluga | 29 |
| 8.4 | Challenge: Benford's law | 31 |
| 9 | Chapter 9: Questions and remarks | 35 |
| 9.1 | Simple questions | 35 |
| 9.2 | BH exercises: hints and solutions | 38 |
| 9.3 | Challenge: Betting | 39 |
| 10 | Chapter 10: Questions and remarks | 41 |
| 10.1 | Simple questions | 41 |
| 10.2 | BH exercises: hints and solutions | 44 |
| 10.3 | Challenge: Records | 44 |
| 11 | Old exam questions | 47 |
| 12 | Hints | 51 |
| 13 | Solutions | 59 |

INTRODUCTION

This study guide contains three parts per chapter of BH.

1. The first section contains simple questions and exercises related to the text per section of BH. They are meant to practice while you read BH. These questions are (often much) simpler than exam questions.
2. There is a part related to the obligatory exercises of BH to provide motivational comments, hints and solutions.
3. The third part contains challenges. These problems are (quite a bit) above exam level, hence optional. However, if you like to be intellectually challenged, then you'll like these problems a lot.

In general, when working on these exercises, try first hard to find the answer. Then *write it down* on paper. And only after having written your answer on paper, meticulously compare your work against ours. Like this you'll get a lot of feedback, and you'll see that it is quite hard to get the details right. You should spend time thinking about the definition and notation we use, and understand any differences. Mind that good notation and good understanding strongly correlate.

The last chapter contains old exam questions.

QUESTIONS AND REMARKS

7.1 SIMPLE QUESTIONS

Section 7.1

Ex 7.1.1. In your own words, explain what is

1. a joint PMF, PDF, CDF;
2. a conditional PMF, PDF, CDF;
3. a marginal PMF, PDF, CDF.

Ex 7.1.2. Suppose a head twice out of two coin flips is $P\{X_1 = H, X_2 = H\}$. What has this to do with joint PMFs? Can you generalize this idea to other examples?

Ex 7.1.3. In the previous exercise, suppose the outcome of the second throw is always equal to that of the first. Specify the joint PMF.

Ex 7.1.4. We have two rvs $X, Y \in [0, 1]^2$ (here $[0, 1]^2 = [0, 1] \times [0, 1]$) with the joint PDF $f_{X,Y}(x, y) = 2I_{x \leq y}$.

1. Are X and Y independent?
2. Compute $F_{X,Y}(x, y)$.

Ex 7.1.5. We have two continuous rvs X, Y . Suppose the joint CDF factors into the product of the marginals, i.e., $F_{X,Y}(x, y) = F_X(x)F_Y(y)$. Can it still be possible in general that the joint PDF does not factor into a product of marginal PDFs of X and Y , i.e., $f_{X,Y}(x, y) \neq f_X(x)f_Y(y)$?

Ex 7.1.6. BH define the conditional CDF given an event A on page 416 as $F(y|A)$. Use this definition to write $F_{X,Y}(x, y)/F_X(x)$ as a conditional CDF. Is this equal to the conditional CDF of Y given X ?

Ex 7.1.7. Let X be uniformly distributed on the set $\{0, 1, 2\}$ and let $Y \sim \text{Bern}(1/4)$; X and Y are independent.

1. Present a contingency table for X and Y .
2. What is the interpretation of the column sums of the table?
3. What is the interpretation of the row sums of the table?

4. Suppose you would change some of the entries in the table. Are X and Y still independent?

Ex 7.1.8. A machine makes items on a day. Some items, independent of the other items, are failed (i.e., do not meet the quality requirements). What are N and p in the context of the chicken-egg story of BH? What are the ‘eggs’ in this context, and what is the meaning of ‘hatching’? What type of ‘hatching’ do we have here?

Ex 7.1.9. We have two rvs X and Y on \mathbb{R}^+ . It is given that $F_{X,Y}(x,y) = F_X(x)F_Y(y)$ for $x,y \leq 1/3$. It is true that then X and Y are necessarily independent.

Ex 7.1.10. I select a random guy from the street, his height $X \sim \text{Norm}(1.8, 0.1)$, and I select a random woman from the street, her height is $Y \sim \text{Norm}(1.7, 0.08)$. I claim that since I selected the man and the woman independently, their heights are independent. Briefly comment on this claim.

Ex 7.1.11. For any two rvs X and Y on \mathbb{R}^+ with marginals F_X and F_Y , can it hold that $P\{X \leq x, Y \leq y\} = F_X(x)F_Y(y)$?

Ex 7.1.12. Theorem 7.1.11. What is the meaning of the notation $X|N = n$?

Ex 7.1.13. Let X, Y be two discrete rvs with CDF $F_{X,Y}$. Can we compute the PDF as $\partial_x \partial_y F_{X,Y}(x,y)$?

Ex 7.1.14. Redo BH.7.1.24 with indicator functions and the fundamental bridge (recall, $P\{A\} = E[I_A]$ for an event A). (Indicators are often easy to use, and prevent many mistakes, as is demonstrated with this example.)

Section 7.2

Ex 7.1.15. BH.7.2.2. Write down the integral to compute $E[(X - Y)^2]$, and solve it.

Ex 7.1.16. Explain that for a continuous r.v. X with CDF F and a and b (so it might be that $a > b$),

$$P\{a < X < b\} = [F(b) - F(a)]^+.$$
 (7.1.1)

Remark 7.1.17. If you are like me, you underestimate at first the importance of using indicator functions. In fact, they are extremely useful for several reasons.

1. They help to keep your formulas clean.
2. You can use them in computer code as logical conditions, or to help counting relevant events, something you need when numerically estimating multi-D integrals, for machine learning for instance.
3. Even though figures give geometrical insight into how to integrate over an 2D area, when it comes to reversing the sequence of integration, indicators are often easier to use.

4. In fact, *expectation is the fundamental concept in probability theory*, and the *probability of an event is defined as*

$$P\{A\} := E[I_A]. \quad (7.1.2)$$

Thus, the fundamental bridge is actually an application of LOTUS to indicator functions. Hence, reread BH.4.4!

Ex 7.1.18. What is $\int_{-\infty}^{\infty} I_{0 \leq x \leq 3} dx$?

Ex 7.1.19. What is

$$\int x I_{0 \leq x \leq 4} dx? \quad (7.1.3)$$

When we do an integral over a 2D surface we can first integrate over the x and then over the y , or the other way around, whatever is the most convenient. (There are conditions about how to handle multi-D integral, but for this course these are irrelevant.)

Ex 7.1.20. What is

$$\iint xy I_{0 \leq x \leq 3} I_{0 \leq y \leq 4} dx dy? \quad (7.1.4)$$

Ex 7.1.21. What is

$$\iint I_{0 \leq x \leq 3} I_{0 \leq y \leq 4} I_{x \leq y} dx dy? \quad (7.1.5)$$

Ex 7.1.22. Take $X \sim \text{Unif}([1, 3])$, $Y \sim \text{Unif}([2, 4])$ and independent. Compute

$$P\{Y \leq 2X\}. \quad (7.1.6)$$

Section 7.3

Ex 7.1.23. Give a brief example of a situation where it might be more convenient to employ the correlation than the covariance. Explain why.

Ex 7.1.24. In queueing theory the concept of squared coefficient of variance SCV of a rv X is very important. It is defined as $C = V[X]/(E[X])^2$. Is the SCV of X equal to $\text{Corr}(X, X)$? Can it happen that $C > 1$?

Ex 7.1.25. Prove the key properties 1–5 of the covariance below BH.7.3.2.

Ex 7.1.26. Using the definition of Covariance (BH.7.3.1) derive the expression $\text{Cov}[X, Y] = E[XY] - E[X]E[Y]$. Use this to show why independence of X and Y implies their uncorrelatedness (Note that the converse does not hold).

Ex 7.1.27. Let U, V be two rvs and let $a, b \in \mathbb{R}$. Using the previous question to express $\text{Cov}[a(U + V), b(U - V)]$ in terms of $V[U]$, $V[V]$ and $\text{Cov}[U, V]$.

Ex 7.1.28. The solution of BH.7.3.6 is a somewhat tricky; I would have not found this trick myself. Here is an approach that is trick free.

Neglecting the event $\{X = Y\}$ as this has zero probability, we know that $M = X, L = Y$ or $M = Y, L = X$. Use this idea and the formula $\text{Cov}[M, L] = E[ML] - E[M]E[L]$ to derive the result of this example.

Section 7.4

Ex 7.1.29. Come up with a short illustrative example in which the random vector $\mathbf{X} = (X_1, \dots, X_6)$ follows a Multinomial Distribution with parameters $n = 10$ and $\mathbf{p} = (\frac{1}{6}, \dots, \frac{1}{6}) \in \mathbb{R}^6$.

Section 7.5

Ex 7.1.30. Is the following claim correct? If the rvs X, Y are both normally distributed, then (X, Y) follows a Bivariate Normal distribution.

Ex 7.1.31. Let X, Y, Z be iid $\mathcal{N}(0, 1)$. Determine whether or not the random vector

$$\mathbf{W} = (X + 2Y, 3X + 4Z, 5Y + 6Z, 2X - 4Y + Z, X - 9Z, 12X + \sqrt{3}Y - \pi Z + 18)$$

is Multivariate Normal. (Explain in words, don't do a lot of tedious math here!)

7.2 MEMORYLESS EXCURSIONS: A CONFUSING PROBLEM WITH MEMORYLESS RVS

BY give a quick argument to compute $E[M]$ and $E[L]$ where $M = \max\{X, Y\}$ and $L = \min\{X, Y\}$ are the maximum and minimum of two iid exponential rvs X and Y . Since X and Y have the same distribution,

$$E[L] + E[M] = E[L + M] = E[X + Y] = 2E[X].$$

Therefore,

$$E[M] = 2E[X] - E[L]. \quad (7.2.1)$$

Next, by the fact that X and Y are memoryless,

$$E[M] = E[L] + E[X]. \quad (7.2.2)$$

An interpretation can help to see this. There are two machines, each working on a job in parallel. Let X and Y be the production times at either machine. The time the first job finishes is evidently $L = \min\{X, Y\}$. Then, *due to memorylessness*, the service time of the remaining job 'restarts'; this takes an expected time $E[X]$ to complete. Adding these two equations and noting that $E[L]$ cancels we get $2E[M] = 3E[X]$, hence:

$$E[M] = \frac{3}{2}E[X], \quad E[L] = E[M] - E[X] = \frac{1}{2}E[X]. \quad (7.2.3)$$

This argument seems general enough, so it must hold for discrete memoryless rvs too, i.e., when $X, Y \sim \text{Geo}(p)$. But that is not the case: it is only true when $X, Y \sim \text{Exp}(\lambda)$ and independent. To see what is wrong I tried as many different approaches to this problem I

could think of, which resulted in this text.¹ In Section 7.2.1 we'll derive (7.2.1) for geometric rvs in multiple different ways. Hence, the culprit must be (7.2.2). Then, in Section 7.2.2 we'll show that both equations *are true* for exponential rvs. Finally, in Section 7.2.3 we find a formula that is similar to $E[M] = E[L] + E[X]$ but that holds for both types of memoryless rvs, whether they are discrete or continuous.

THE ANALYSIS OF the above problem illustrates many general and useful probability concepts such as joint CDF, joint PMF/PDF, the fundamental bridge, integration over 2D areas, 2D LOTUS, conditional PMF/PDF, MGFs, and the change of variables formula. It pays off to do the exercises yourself and then study the hinta and solutions carefully.

YOU'LL NOTICE, hopefully, that I use many different methods to the same problem, and that I take pains to see how the answers of these methods relate. There are at least two reasons for this. Often, a problem can be solved in multiple ways, and one method is not necessarily better than another; better yet, different methods may augment the understanding of the problem. The second reason is that it is easy to make a mistake in probability. If different methods give the same answer, the probability of having made a mistake becomes smaller.

7.2.1 Discrete memoryless rvs

Before embarking on a problem, it often helps to refresh our memory. This is what we do first. Let X be $\sim \text{Geo}(p)$ and write $q = 1 - p$.

Ex 7.2.1. What is the domain of X ?

With some fun tricks with recursions it is possible to quickly derive the most important expressions for geometric rvs:

$$P\{X > 0\} = P\{\text{failure}\} = q$$

$$P\{X > j\} = q P\{X > j-1\} \implies P\{X > j\} = q^j P\{X > 0\} = q^{j+1}.$$

$$P\{X \geq j\} = P\{X > j-1\} = q^j.$$

$$P\{X = j\} = P\{X > j-1\} - P\{X > j\} = q^j - q^{j+1} = (1-q)q^j = pq^j.$$

$$E[X] = p \cdot 0 + q(1 + E[X]) \implies E[X] = q/(1-q) = q/p.$$

Mind that, even though this is neat, it only work for geometric rvs.

Ex 7.2.2. Explain the above.

Clearly, such tricks are nice and quick, but they are not general. We should also practice with the general method.

¹ Part of this material was born out of annoyance. The book uses one of those typical probability arguments: slick, half complete, and wrong as soon as one tries it in other situations. In other words, the type of argument beginner books should stay clear of. I admit that I was quite irritated about the argument offered by the book.

Ex 7.2.3. Simplify $P\{X > j\} = \sum_{i=j+1}^{\infty} P\{X = i\}$ to see that this is equal to q^{j+1} . Realize that from this, $P\{X \geq j\} = P\{X > j-1\} = q^j$.

Ex 7.2.4. Use indicator variables show that

$$E[X] = \sum_{i=0}^{\infty} i P\{X = i\} = \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} I_{j < i} P\{X = i\} = p/q.$$

Ex 7.2.5. Look up the definition of a memoryless rv, and check that X is memoryless.

WITH THIS REFRESHER, we can derive some useful properties of the minimum $L = \min\{X, Y\}$, where $Y \sim \text{Geo}(p)$ and independent of X . For this we use the fundamental bridge and 2D LOTUS, which in general read like

$$P\{g(X, Y) \in A\} = E[I_{g(X, Y) \in A}] = \sum_i \sum_j I_{g(i, j) \in A} P\{X = i, Y = j\}.$$

Ex 7.2.6. What is the domain of L ? Then, show that

$$P\{L \geq i\} = q^{2i} \implies L \sim \text{Geo}(1 - q^2).$$

Ex 7.2.7. Show that

$$E[L] = q^2/(1 - q^2).$$

Ex 7.2.8. Show that

$$E[L] + E[X] = \frac{q}{1 - q} \frac{1 + 2q}{1 + q}. \quad (7.2.4)$$

NOW WE CAN combine these facts with the properties of the maximum $M = \max\{X, Y\}$.

Ex 7.2.9. Show that

$$2E[X] - E[L] = \frac{q}{1 - q} \frac{2 + q}{1 + q}.$$

Clearly, unless $q = 0$, $E[L] + E[X] \neq 2E[X] - E[L]$, hence, $E[M]$ can only be one of the two and (7.2.1) and (7.2.2) cannot be both true.

To convince ourselves that [7.2.9], hence (7.2.1), is indeed true, we pursue three ideas.

HERE IS THE FIRST idea.

Ex 7.2.10. Show for the PMF of M that

$$p_M(k) = P\{M = k\} = 2pq^k(1 - q^k) + p^2q^{2k}.$$

Ex 7.2.11. With the previous exercise, show now that $p_M(k) = 2P\{X = k\} - P\{L = k\}$.

Ex 7.2.12. Finally, show that $E[M] = 2E[X] - E[L]$.

THE SECOND IDEA.

Ex 7.2.13. First show that $P\{M \leq k\} = (1 - q^{k+1})^2$.

Ex 7.2.14. Simplify $P\{M = k\} = P\{M \leq k\} - P\{M \leq k - 1\}$ to see that $p_M(k) = 2P\{X = k\} - P\{L = k\}$.

AND HERE IS the third idea.

Ex 7.2.15. Explain that

$$P\{L = i, M = k\} = 2p^2 q^{i+k} I_{k>i} + p^2 q^{2i} I_{i=k}.$$

Ex 7.2.16. Use [7.2.15] and marginalization to compute the marginal PMF $P\{M = k\}$.

Ex 7.2.17. Use [7.2.15] to compute $P\{L = i\}$.

IN CONCLUSION, we verified the correctness of $E[M] = 2E[X] - E[L]$ in three different, and useful ways. Let us now focus on exponential rvs rather than geometric rvs.

7.2.2 Continuous memoryless rvs

In this section we analyze the correctness of (7.2.1) and (7.2.2) for continuous memoryless rvs, i.e., exponentially distributed rvs. I decided to analyze this in as much detail as I could think of, hoping that this would provide me with a lead to see how to generalize the equation $E[M] = E[L] + E[X]$ such that it covers also the case with geometric rvs.

FIRST WE NEED to recall some basic facts about the exponential distribution.

Ex 7.2.18. Show that X is memoryless.

Ex 7.2.19. Show that $E[X] = 1/\lambda$.

NOW WE CAN shift our attention to the rvs L and M .

Ex 7.2.20. Show that $F_L(x) = 1 - e^{-2\lambda x}$.

Clearly, this implies that $L \sim \text{Exp}(2\lambda)$ and $E[L] = 1/(2\lambda) = E[X]/2$. Hence, we see that (7.2.3) holds now. Moreover, with the same trick we see that the distribution function F_M for the maximum M is given by

$$F_M(v) = P\{M \leq v\} = P\{X \leq v, Y \leq v\} = (F_X(v))^2 = (1 - e^{-\lambda v})^2.$$

Of course this is a nice trick, but it is not a method that allows us to compute the distribution for more general functions of X and Y . For more general cases, we have to use the fundamental bridge and LOTUS, that is, for any set² A in the domain of $X \times Y$

$$\begin{aligned} P\{g(X, Y) \in A\} &= E[I_{g(X, Y) \in A}] = \iint I_{g(x, y) \in A} f_{X, Y}(x, y) dx dy \\ &= \iint_{g^{-1}(A)} f_{X, Y}(x, y) dx dy = \iint_{\{(x, y): g(x, y) \in A\}} f_{X, Y}(x, y) dx dy. \end{aligned}$$

² If you like maths, you should be quite a bit more careful about what type of set A is acceptable. Here such matters are of no importance.

The joint CDF $F_{X,Y}$ then follows because $F_{X,Y}(x,y) = P\{X \leq x, Y \leq y\} = E[I_{X \leq x, Y \leq y}]$. A warning is in place: conceptually this approach is easy, but doing the integration can be very challenging (or impossible). However, this expression is very important as this is the preferred way to compute distributions by numerical methods and simulation.

Ex 7.2.21. Use the fundamental bridge to rederive the above expression for $F_M(v)$.

Ex 7.2.22. Show that the density of M has the form $f_M(v) = \partial_v F_M(v) = 2(1 - e^{-\lambda v})\lambda e^{-\lambda v}$.³

Ex 7.2.23. Use the density f_M to show that $E[M] = 2E[X] - E[L]$.

Recalling that we already obtained $E[L] = E[X]/2$, we see that $E[M] = 2E[X] - E[L] = 3E[X]/2$, which settles the truth of (7.2.3).

WE CAN ALSO compute the densities $f_M(y)$ (and $f_L(x)$ by marginalizing the joint density $f_{L,M}(x,y)$. However, for this, we first need the joint distribution $F_{L,M}$, and then we can get $f_{L,M}$ by differentiation, i.e., $f_{X,Y} = \partial_x \partial_y F_{X,Y}$. Let us try this approach too.

Ex 7.2.24. Use the fundamental bridge to show that for $u \leq v$,

$$F_{L,M}(u,v) = P\{L \leq u, M \leq v\} = 2 \int_0^u (F_Y(v) - F_Y(x))f_X(x)dx.$$

Ex 7.2.25. Take partial derivatives to show that

$$f_{L,M}(u,v) = 2f_X(u)f_Y(v)I_{u \leq v}.$$

Ex 7.2.26. In [7.2.25] marginalize out L to find f_M , and marginalize out M to find f_L .

WE DID A number of checks for the case $X, Y, \text{iid}, \sim \text{Exp}(\lambda)$, but I have a second way to check the consistency of our results. For this I use the idea that the geometric distribution is the discrete analog of the exponential distribution. Now we study how this works, and that by taking proper limits we can obtain the results for the continuous setting from the discrete setting.

First, let's try to obtain an intuitive understanding of how $X \sim \text{Geo}(\lambda/n)$ approaches $Y \sim \text{Exp}(\lambda)$ as $n \rightarrow \infty$. For this, divide the interval $[0, \infty)$ into many small intervals of length $1/n$. Let $X \sim \text{Geo}(\lambda/n)$ for some $\lambda > 0$ and $n \gg 0$. Then take some $x \geq 0$ and let i be such that $x \in [i/n, (i+1)/n)$.

Ex 7.2.27. Show that

$$P\{X/n \approx x\} \approx \frac{\lambda}{n} \left(1 - \frac{\lambda}{n}\right)^{xn}. \quad (7.2.5)$$

Next, introduce the highly intuitive notation⁴ $dx = \lim_{n \rightarrow \infty} 1/n$, and use the standard limit⁵ $(1 - \lambda/n)^n \rightarrow e^{-\lambda}$ as $n \rightarrow \infty$ to see that (7.2.5) converges to

$$P\{X/n \approx x\} \rightarrow \lambda e^{-\lambda x} dx = f_X(x)dx, \quad \text{as } n \rightarrow \infty.$$

³ We write ∂_v as a shorthand for d/dv in the 1D case, and for ∂/∂_v the partial derivative in the 2D case.

⁴ In your math classes you learned that $\lim_{n \rightarrow \infty} 1/n = 0$. Doesn't this definition therefore imply that $dx = 0$? Well, no, because dx is not a real number but an infinitesimal. Infinitesimals allow us to consider a quantity that is so small that it cannot be distinguished from 0 within the real numbers.

⁵ This is not entirely trivial to prove. If you like mathematics, check the neat proof in Rudin's Principles of mathematical analysis.

If you don't like this trick with dx , here is another method, based on with moment-generating functions.

Ex 7.2.28. Derive the moment-generating function $M_{X/n}(s)$ of X/n when $X \sim \text{Geo}(p)$. Then, let $p = \lambda/n$, and show that $\lim_{n \rightarrow \infty} M_{X/n}(s) = M_Y(s)$, where $Y \sim \text{Exp}(\lambda)$.

With these limits in place, we can relate the minimum $L = \min\{X, Y\}$ for the discrete and the continuous settings.

Ex 7.2.29. Suppose that $X, Y \sim \text{Geo}(\lambda/n)$, then check that $\lim_{n \rightarrow \infty} E[L/n] = 1/2\lambda$.

Clearly, $1/2\lambda = E[X]/2$ when $X \sim \text{Exp}(\lambda)$.

Here is yet another check on the correctness of $f_M(x)$.

Ex 7.2.30. Show that the PMF $P\{M = k\}$ for the discrete M in [7.2.10] converges to $f_M(x)$ of [7.2.22] when $n \rightarrow \infty$. Take k suitable.

FINALLY, I HAVE a third way to check the above results, namely by verifying (7.2.2), i.e. $E[M - L] = E[X]$. For this, we compute the joint CDF $F_{L, M-L}(x, y)$. With this, you'll see directly how to compute $E[M - L]$.

Ex 7.2.31. Use the fundamental bridge to obtain

$$F_{L, M-L}(x, y) = (1 - e^{-2\lambda x})(1 - e^{-\lambda y}) = F_L(x)F_Y(y).$$

Ex 7.2.32. Conclude that $M - L$ and L are independent, and $M - L \sim Y$.

By the above exercise, we find that $E[M - L] = E[Y] = E[X]$, as X and Y are iid.

THIS MAKES ME wonder whether $M - L$ and L are also independent for the discrete case, i.e., when X, Y iid and $\sim \text{Geo}(p)$. Hence, we should check that for all i, j

$$P\{L = i, M - L = j\} = P\{L = i\} P\{M - L = j\}. \quad (7.2.6)$$

Ex 7.2.33. Use [7.2.15] to see that

$$P\{L = i, M - L = j\} = 2p^2 q^{2i+j} I_{j \geq 1} + p^2 q^{2i} I_{j=0}.$$

Now for the RHS.

Ex 7.2.34. Derive that

$$P\{M - L = j\} = \frac{2p^2 q^j}{1 - q^2} I_{j \geq 1} + \frac{p^2 q^j}{1 - q^2} I_{j=0}.$$

Recalling that $P\{L = i\} = (1 - q^2)q^{2i}$, it follows right away from (7.2.6) that L and $M - L$ are independent. Interestingly, from [7.2.31] we see $M - L \sim Y$ for the continuous case. However, here, for the geometric case, $P\{M - L = j\} \neq p q^j = P\{Y = j\}$. This explains why $E[M] \neq E[L] + E[X]$ for geometric rvs: we should be more carefull in how to split M in terms of L and X .

ALL IN ALL, we have checked and double checked all our expressions and limits for the geometric and exponential distribution. We had success too: the solution of the last exercise provides the key to understand why (7.2.1) and (7.2.2) are true for exponentially distributed rvs, but not for geometric random variables. In fact, in the solutions we can see the term corresponding to $X = Y = i$ for $X, Y \sim \text{Geo}(p)$ becomes negligibly small when $n \rightarrow 0$. In other words, $P\{X = Y\} > 0$ when X and Y are discrete, but $P\{X = Y\} = 0$ when X and Y are continuous. Moreover, by [7.2.31], $E[M] = E[L + M - L] = E[L] + E[M - L]$, but $E[M - L] \neq E[X]$. So, to resolve our leading problem we should reconsider $E[M - L]$.

7.2.3 The solution

Let us now try to repair (7.2.2), i.e., $E[M] = E[L] + E[X]$, for the case $X, Y \sim \text{Geo}(p)$. We should be careful about the non-negligible case that $M = L$, so we move, carefully, step by step.

Ex 7.2.35. Why is the following true:

$$E[M] = E[L] + E[(M - L)I_{M > L}] = E[L] + 2E[(Y - X)I_{Y > X}]. \quad (7.2.7)$$

Ex 7.2.36. Show that

$$2E[(Y - X)I_{Y > X}] = \frac{2q}{1 - q^2}.$$

Ex 7.2.37. Combine the above with the expression for $E[L]$ of [7.2.7] to obtain [7.2.9] for $E[M] = 2E[X] - E[L]$, thereby verifying the correctness of (7.2.7).

While (7.2.7) is correct, I am still not happy with the second part of (7.2.7) as I find it hard/unintuitive to interpret. Restarting again from scratch, here is another attempt to rewrite $E[M]$ by using $Z \sim \text{FS}(p)$, i.e., Z has the first success distribution with parameter p , in other words, $Z \sim X + 1$ with $X \sim \text{Geo}(p)$.

Ex 7.2.38. Explain that

$$E[M] = E[L] + E[Z I_{M > L}], \quad (7.2.8)$$

Ex 7.2.39. Show that

$$E[Z I_{M > L}] = \frac{2q}{1 - q^2},$$

i.e., the same as [7.2.36], hence (7.2.8) is correct.

I AM NEARLY happy, but I want to see that (7.2.8), which is correct for discrete rvs, has also the correct limiting behavior.

Ex 7.2.40. Show that $E[Z/n I_{M > L}] \rightarrow 1/\lambda$, which is the expectation of an $\text{Exp}(\lambda)$ rv!

Finally I understand why $E[M] = E[L] + E[X]$ for $X, Y \sim \text{Exp}(\lambda)$ but not for when X, Y are discrete. For discrete rvs, L and M can be equal, while for continuous rvs, this

is impossible⁶ It took a long time, and a lot of work, to understand how to resolve the confusing problem, but I learned a lot. In particular, I find (7.2.8) a nice and revealing equation.

Finally, if you like to train with the tools you learned, you can try your hand at analyzing the same problem, but now for uniform X, Y .

7.3 EXERCISES ON 2D INTEGRATION

Here is some extra material for you practice on 2D integration, indicators and 2D LOTUS. These exercises are old exam questions, hence quite a bit harder than the above. They form important training.

Ex 7.3.1. Let X and Y be continuous random variables. Furthermore, $F(x, y)$ is the joint cumulative distribution function of X and Y . This function has the following properties.

1. $F(x, y) = \frac{1}{8}(x-1)^2(y-2)$ for $1 < x < 3$ and $2 < y < 4$,
2. $\frac{\partial F(x, y)}{\partial x} = 0$ for $x \notin (1, 3)$,
3. $\frac{\partial F(x, y)}{\partial y} = 0$ for $y \notin (2, 4)$.

Use these properties to answer the following questions.

1. What is $F(2, 5)$?
2. Determine the joint probability density function of X and Y .
3. Determine $P(2 < X < 3, 2 < Y < 4)$.
4. Determine the joint probability $P(Y < 2X, 2X + Y > 6)$. Clearly draw the area over which you integrate.

Ex 7.3.2. Suppose X and V are independent, $X \sim \text{Expo}(\lambda)$ and $V \sim \text{Expo}(\mu)$. Define the ratio $R = X/V$ and derive the cumulative distribution function (CDF) of R . Provide at least two checks on the CDF to make sure that your result is indeed a valid CDF. Note: There is no need to derive the probability density function (PDF) of R .

Ex 7.3.3. Consider the following joint density function

$$f_{X,Y}(x, y) = \begin{cases} cxy & \text{for } 0 \leq x < \frac{1}{2} \text{ and } 0 \leq y \leq x, \\ cxy & \text{for } \frac{1}{2} \leq x \leq 1 \text{ and } 0 \leq y \leq 1-x, \\ 0 & \text{otherwise.} \end{cases}$$

1. What is the correct value of the constant c ?

⁶ A bit more carefully formulated: the event $\{L = M\}$ has zero probability for continuous rvs.

2. Derive the conditional probability density function $f_{X|Y}(x|y)$. Verify that your result is indeed a valid density function.

Ex 7.3.4. Suppose the random variables X and Y have the joint probability density function

$$f_{X,Y}(x,y) = \begin{cases} 2 & \text{if } 0 \leq x < 1, \quad |y| < \frac{1}{2}(1-x), \\ 0 & \text{otherwise} \end{cases}$$

1. Calculate the marginal probability density functions $f_X(x)$ and $f_Y(y)$ and show that these are valid probability density functions.
2. Find the conditional expectation $E[X|Y = y]$. Provide at least one ‘sanity check’ that shows that your answer makes intuitive sense. If you did not find an answer to (a), you can use that

$$f_Y(y) = \begin{cases} 2(1-2|y|) & \text{if } -\frac{1}{2} < y < \frac{1}{2}, \\ 0 & \text{otherwise.} \end{cases}$$

3. Calculate the joint cumulative distribution function $F_{X,Y}(x,y)$ for $x = \frac{1}{2}$ and $y = 3$.

Ex 7.3.5. Suppose the joint probability density function of X and Y is given by

$$f_{X,Y}(x,y) = \frac{c}{1-x}, \quad 0 < x+y < 1, \quad x > 0, \quad y > 0$$

and $f_{X,Y}(x,y) = 0$ otherwise.

1. For what value of the constant c is $f_{X,Y}(x,y)$ a joint probability density function?
2. What is the probability that $X + Y > \frac{1}{2}$?
3. What is the probability that both X and Y are smaller than $\frac{1}{2}$ given that $X + Y > \frac{1}{2}$?

Ex 7.3.6. X and Y be random variables with joint probability density function

$$f_{X,Y}(x,y) = \begin{cases} \frac{3}{16}xy^2, & 0 \leq x \leq c \text{ and } 0 \leq y \leq c, \\ 0, & \text{elsewhere} \end{cases}$$

where $c > 0$ is a real number.

1. Show that $c = 2$.
2. Show that $P(X + Y > 2) = \frac{9}{10}$. Start by making a clear sketch of the area in the (x,y) -plane over which you take the required integral.
3. Calculate the conditional probability $P(Y < X^2 | X + Y < 2)$.

7.4 BH EXERCISES: HINTS AND SOLUTIONS

Ex 7.4.1. BH.7.1. We simulate this in one of the assignments.

Ex 7.4.2. BH 7.9. We'll develop a simulation for this in the assignments.

Ex 7.4.3. BH.7.10.

Recall that a conditional CDF given an event A is defined as $F(y|A) = P\{Y \leq y|A\}$. Likewise, let us write here $F_T(t|x) = P\{T \leq t|X = x\}$. Just use this in your derivation. However, there is one problem with the fact that the event $\{X = x\}$ has probability zero. In the solution I'll discuss how to around this.

Don't forget to compare this exercise to BH.7.9, which is the same but for discrete memoryless rvs.

Ex 7.4.4. BH.7.11

Ex 7.4.5. BH.7.13

Ex 7.4.6. BH.7.15.

Ex 7.4.7. BH.7.24. In the assignments we'll develop a simulator.

Ex 7.4.8. BH.7.29

Ex 7.4.9. BH. 7.38. Besides the solution of BH, read our solution.

Ex 7.4.10. BH.7.53. We simulate this in one of the assignments. The ideas of this exercise find much use in finance, physics, and actuarial sciences. In particular, the expected time it takes the drunken person—It's not only guys that sometimes consume too much alcohol—to hit some boundary is interesting. The notation of the book is a bit clumsy. Here is better notation. Let X_i be the movement along the x -axis at step i , and Y_i along the y -axis. Then $S_n = \sum_{i=1}^n X_i$ and $T_n = \sum_{j=1}^n Y_j$, and $R_n^2 = S_n^2 + T_n^2$.

Ex 7.4.11. BH.7.58. This is a totally great exercise. First solve it yourself. In the solution, I'll explain why, in particular how to relate the concept of covariance to the determinant of a matrix.

Ex 7.4.12. BH.7.59. Read this exercise, then read (and do) BH.5.53 for some further background. You'll encounter these topics countless times in other courses! The final answer is really nice and intuitive.

Ex 7.4.13. BH.7.71.

Ex 7.4.14. BH.7.86. The concepts discussed here are a standard part of the education of GPs (i.e., medical doctors), and in data science in general.

7.5 CHALLENGE: A UNIQUENESS PROPERTY OF THE POISSON DISTRIBUTION

Consider the chicken-egg story (BH 7.1.9): A chicken lays a random number of eggs N and each egg independently hatches with probability p and fails to hatch with probability $q = 1 - p$. Formally, $X|N \sim \text{Bin}(N, p)$. Assume also that $X|N \sim \text{Bin}(N, p)$ and that $N - X$ is independent of X . For $N \sim \text{Pois}(\lambda)$ it is shown in BH 7.1.9 that X and Y are independent. This exercise asks for the converse: showing that the independence of X and Y implies that $N \sim \text{Pois}(\lambda)$ for some λ . Hence, the Poisson distribution is quite special: it is the only distribution for which the number of hatched eggs doesn't tell you anything about the number of unhatched eggs.

Let $0 < p < 1$. Let N be an rv. taking non-negative integer values with $P(N > 0) > 0$. Assume also that $X|N \sim \text{Bin}(N, p)$ and that $N - X$ is independent of X .

Ex 7.5.1. Use the assumption that $P\{N > 0\} > 0$ to prove that N has support \mathbb{N} , i.e. $P\{N = n\} > 0$ for all $n \in \mathbb{N}$. Note: $0 \in \mathbb{N}$.

Ex 7.5.2. Write $Y = N - X$. Prove that

$$P\{X = x\} P\{Y = y\} = \binom{x+y}{x} p^x (1-p)^y P\{N = x+y\}. \quad (7.5.1)$$

Ex 7.5.3. Prove that N is Poisson distributed.

7.6 CHALLENGE: IMPROPER INTEGRALS AND THE CAUCHY DISTRIBUTION

This problem challenges your integration skills and lets you think about the subtleties of integrating a function over an infinite domain.⁷

Assume that X has the Cauchy distribution. Recall that $E[X]$ does not exist (hence, it is not automatic that the expectation of a some arbitrary rv. exists).

Ex 7.6.1. Why does $E\left[\frac{|X|}{X^2+1}\right]$ exist? Find its value. It is essential that you include your arguments.

Ex 7.6.2. Explain why the previous exercise implies that $E\left[\frac{X}{X^2+1}\right]$ exists. Then find its value.

7.7 CHALLENGE: PROOF ABOUT INDEPENDENCE OF NORMAL RVS

Consider two iid rvs X, Y such that $X + Y$ and $X - Y$ are independent. In BH 7.5.8, it is claimed that this implies that X and Y are normally distributed.

This challenge asks to give a proof of this claim. Throughout this problem, you may assume that X and Y have a MGF that is defined for all $t \in \mathbb{R}$.⁸ You may also use without

⁷ Such integrals are called improper Riemann integrals.

⁸ This may seem like a big restriction, but this argument can easily be adapted to work with the *characteristic function* instead of the MGF, and the characteristic function does always exist. You will learn the characteristic function in the second year courses Statistical Inference and Linear Models in Statistics.

proof the fact that MGFs that are defined everywhere are infinitely often continuously differentiable.

Ex 7.7.1. Let M_X be the MGF of X (and hence of Y).

1. Prove that $M_X(2t) = (M_X(t))^3(M_X(-t))$.
2. Define $f(t) = \log M_X(t)$. Prove that $8f'''(2t) = 3f'''(t) - f'''(-t)$.
3. Let $R > 0$ be arbitrary. Use Weierstrass' theorem to prove that f''' attains a minimum m and a maximum M on the interval $[-R, R]$, and then prove that $m = M = 0$.
4. Prove that X is normally distributed.

7.8 CHALLENGE: RECOURSE MODELS

This exercise will give an example of how probability theory can pop up in OR problems, in particular in linear programs. It introduces you to the concept of *recourse models*, which you will learn about in the master course Optimization Under Uncertainty. Disclaimer: the story is quite lengthy, but the concepts introduced and questions asked are in fact not very hard. We just added the story to make things more intuitive.

WE CONSIDER A pastry shop that only sells one product: chocolate muffins. Every morning at 5:00 a.m., the shop owner bakes a stock of fresh muffins, which he sells during the rest of the day. Making one muffin comes at a cost of $c = \$1$ per unit. Any leftover muffins must be discarded at the end of the day, so every morning he starts with an empty stock of muffins.

The owner has one question for you: determine the amount x of muffins that he should make in the morning to minimize his production cost. Note that the owner never wants to disappoint any customer, i.e., he requires that $x \geq d$, where d is the daily demand for muffins.

The problem can be formulated as a linear program (LP):

$$\min_{x \geq 0} \{cx : x \geq d\}. \quad (7.8.1)$$

For simplicity, we ignore the fact that x should be integer-valued.

Ex 7.8.1. Determine the optimal value x^* for x and the corresponding objective value in case d is deterministic.

Of course, in practice d is not deterministic. Instead, d is a random variable with some distribution. However, note that the LP above is ill-defined if d is a random variable. We cannot guarantee that $x \geq d$ if we do not know the value of d .

You explained the issue to the shop owner and he replies: “Of course, you’re right! You know, whenever I’ve run out of muffins and a customer asks for one, I make one on the spot. I never disappoint a customer, you know! It does cost me 50% more money to produce them on the spot, though, you know.”

Mathematically speaking, the shop owner just gave you all the (mathematical) ingredients to build a so-called *recourse model*. We introduce a *recourse variable* y in our model, representing the amount of muffins produced on the spot. Production comes at a unit cost of $q = 1.5c = \$1.5$. Assuming that we know the distribution of d , we can then minimize the *expected total cost*:

$$\min_{x \geq 0} \{cx + E[v(d, x)]\}, \quad (7.8.2)$$

where $v(d, x)$ is the optimal value of another LP, namely the *recourse problem*:

$$v(d, x) := \min_{y \geq 0} \{qy : x + y \geq d\}, \quad (7.8.3)$$

for given values of d and x . The recourse problem can easily be solved explicitly: we get $y = d - x$ if $d \geq x$ and $y = 0$ if $d < x$. So we obtain

$$v(d, x) = q(d - x)^+, \quad (7.8.4)$$

where the operator $(\cdot)^+$ represents the *positive value operator*, defined as

$$(s)^+ = \begin{cases} s & \text{if } s \geq 0, \\ 0 & \text{if } s < 0. \end{cases} \quad (7.8.5)$$

Ex 7.8.2. To get some more insight into the model, suppose (for now) that $d \sim U\{10, 20\}$. Solve the model, i.e., find the optimal amount x^* .

Ex 7.8.3. What is the expected recourse cost (expected cost of on-the-spot production) at the optimal solution x^* , i.e., compute $E[v(d, x^*)]$?

To solve the model correctly, we need the true distribution of d . We learn the following from the shop owner: “My granddaughter, who’s always running around in my shop, is a bit data-crazy, you know, so she’s been collecting some data. I remember her saying that ‘the demand from male and female customers are both approximately normally distributed, with mean values both equal to 10 and standard deviations of 5’. She also mentioned something about correlation, but I don’t remember exactly, you know. It was either almost 1 or almost -1 . I hope this helps!”

Mathematically, we’ve learned that $d = d_m + d_f$, with $(d_m, d_f) \sim \mathcal{N}(\mu, \Sigma)$, where $\mu = (\mu_m, \mu_f) = (10, 10)$ and $\Sigma_{11} = \sigma_m^2 = \Sigma_{22} = \sigma_f^2 = 5^2 = 25$. Finally, $\Sigma_{12} = \Sigma_{21} = \text{Cov}[d_m, d_f] = \rho\sigma_m\sigma_f = 25\rho$. Also, we know that either $\rho \approx 1$ or $\rho \approx -1$.

Ex 7.8.4. Calculate x^* and the corresponding objective value for the case $\rho = -1$. (Do not read $\rho = 1$, this case is not simple.)

Ex 7.8.5. Consider the two extreme cases $\rho = 1$ and $\rho = -1$. In which case will the shop owner have lower expected total costs? Provide a short, intuitive explanation.

CHAPTER 8: QUESTIONS AND REMARKS

8.1 SIMPLE QUESTIONS

Section 8.1

Ex 8.1.1. In probability theory we often want to study properties of functions of rvs. Provide an example for such a function.

Ex 8.1.2. Let the rv X be uniform on the set $\{0, 1, 2, 3, 4, 5\}$. Derive the PMF and the CDF of $Z = 3X$. Explicitly specify the domain.

Ex 8.1.3. Suppose $y = g(x)$ for some differentiable function g . We like to express the PDF f_Y for $Y = g(X)$ in terms of the PDF f_X and g . This is easy when g is strictly increasing and has an inverse at y , because

$$F_Y(y) = P\{Y \leq y\} = P\{g(X) \leq y\} = P\{X \leq g^{-1}(y)\} = F_X(g^{-1}(y)). \quad (8.1.1)$$

Now we take the derivative at the LHS and RHS to get with the chain rule

$$f_Y(y) = F'_Y(y) = \frac{d}{dy} F_X(g^{-1}(y)) = f_X(g^{-1}(y)) \frac{d}{dy} g^{-1}(y) = f_X(x) \frac{1}{g'(x)},$$

where we write $x = g^{-1}(y)$ in the last step. But why is the derivative of $g^{-1}(y)$ at y equal to $1/g'(x)$, with $x = g^{-1}(y)$?

Ex 8.1.4. When g is not strictly increasing everywhere, there can be no or multiple points x such that $g(x) = y$. Explain that in such cases it is much more difficult to express F_Y in terms of F_X than directly use the densities (assuming that g is differentiable). Extend your reasoning to 2D.

Ex 8.1.5. The general 1D change of variables formula is like this,

$$f_Y(y) = \sum_{x_i: g(x_i)=y} f_X(x_i) \frac{1}{|g'(x_i)|},$$

with some natural conditions on g . Apply this formula to the case $g(x) = x^2$.

Ex 8.1.6. If $X \sim \text{Exp}(1)$, use the change-of-variables theorem to obtain the density of $Y = g(X) = \lambda X$. What is $E[Y]$?

Ex 8.1.7. Show that the 1D change-of-variables formula relates directly to the substitution rule of integration theory to solve 1D integrals.

Ex 8.1.8. Use the change of variable formula to relate the $\text{Geo}(p)$ and the $\text{FS}(p)$ distributions.

Ex 8.1.9. BH.8.1.1 write that ‘The support of Y is all $g(x)$ with x in the support of X .’ Do they say that $\text{supp}(Y) = \{x : g(x) \in \text{supp}(X)\}$? BTW, with is the difference between $\text{supp}(X)$ and $\text{sup } X$?

Ex 8.1.10. BH.8.1.3. Check how all moments were found.

Ex 8.1.11. Let $X \sim \text{Unif}(0, 5)$. Using the one dimensional change of variables theorem (BH.8.1.1), derive the PDF and the CDF of $Z = 3X$. Explicitly specify the domain.

Ex 8.1.12. When $Z = X^3$ and $X \sim \text{Unif}(0, 5)$, using the one dimensional change of variables theorem to derive the PDF and the CDF of Z . Specify the domain of Z .

Ex 8.1.13. Let $X \sim \text{Norm}(\mu, \sigma^2)$. Using the one dimensional change of variables theorem BH.8.1.1, show that $Z = \frac{X - \mu}{\sigma} \sim \text{Norm}(0, 1)$.

Ex 8.1.14. Let $X \sim \text{Exp}(1)$. Derive the PDF of e^{-X} .

Ex 8.1.15. Let X, Y be iid standard normal. Using the n -dimensional change of variables theorem, derive the joint PDF of $(X + Y, X - Y)$.

Check your final answer using BH.7.5.8.

Ex 8.1.16. Specify the domain of the new random variable for the following transformations; this important aspect of the change of variables is often overlooked. Let U, V, W, X, X_1, X_2, Y and Z be rvs and let a, b and c be arbitrary constants.

1. $Z = Y^4$ for $Y \in (-\infty, \infty)$;
2. $Y = X^3 + a$ for $X \in (0, 1)$;
3. $U = |V| + b$ for $V \in (-\infty, \infty)$;
4. $Y = e^{X^3}$ for $X \in (-\infty, \infty)$;
5. $V = U I_{U \leq c}$ for $U \in (-\infty, \infty)$;
6. $Y = \sin(X)$ for $X \in (-\infty, \infty)$;
7. $Y = \frac{X_1}{X_1 + X_2}$ for $X_1 \in (0, \infty)$ and $X_2 \in (0, \infty)$;
8. $Z = \log(UV)$ for $U \in (0, \infty)$ and $V \in (0, \infty)$.

Ex 8.1.17. When adding a different equality, we need to be careful to not create a functional relationship between our two new variables U, V , for example $U = X + Y$ and $V = \sin(X + Y)$, or $U = \frac{X}{Y}$ and $V = \frac{Y}{X}$ for conforming X, Y . What would happen to the determinant of the Jacobian matrix if we did? Why would this happen? Explain in your own words.

Section 8.2

Ex 8.1.18. To find the distribution of a convolution through the change of variables formula, we seem to need to add a ‘redundant’ equality? But why is that? What would be the problem if we do not add this? Explain in your own words.

Ex 8.1.19. In this exercise, we combine what we learned in BH.8.1.4 and BH.8.1.9. Let S be the sum of two iid chi-square distributed variables (with one degree of freedom). Using just these two examples, show that $S \sim \text{Exp}(1/2)$.

Ex 8.1.20. A student has obtained an iid random sample of size 2 from a Cauchy distribution. Let the rvs X and Y model the values of the first and second sample. Since s/he does not know what the mean of a Cauchy distribution is, s/he wants to average the sample to obtain what she thinks is a good estimate of the true mean.

To find the distribution of this sample mean, we need to find an expression for $f_W(w)$, where $W = \frac{X+Y}{2}$.

1. Find an expression for $f_W(w)$ in the form of an integral, but do not solve it.
2. It turns out that if we solve the integral, we get that $f_W(w) = f_X(w)$. The distribution of our sample mean is still Cauchy; we did not obtain a better estimate of the Cauchy mean by calculating the sample mean!

Explain (in your own words) why this makes sense.

Section 8.3

Ex 8.1.21. If $a = b = 1$, why is $\text{Beta}(a, b) = \text{Unif}([0, 1])$?

Ex 8.1.22. If $a = b$, why is $\text{Beta}(a, b)$ symmetric?

Ex 8.1.23. If $a > b$ and $X \sim \text{Beta}(a, b)$, is $E[X] > 1/2$?

Ex 8.1.24. BH.8.3.2, last equation. How do the authors get to the equation

$$\beta(a, b) = \frac{1}{(a+b-1) \binom{a+b-2}{a-1}}?$$

Ex 8.1.25. BH.8.3.3. The authors write that X is not marginally Binomial, but is conditionally Binomial. What is the difference?

Ex 8.1.26. BH.8.3.3. The authors use a smart trick to find an expression for the posterior distribution $f_{p|X=k}$ of p . Use this posterior to derive an expression for $P\{X = k\}$ by using the fact that

$$P\{X = k\} = \frac{\binom{n}{k} p^k (1-p)^{n-k} \text{Beta}(a, b)}{\text{Beta}(a+k, b+n-k)},$$

and simplifying the RHS.

Ex 8.1.27. BH.8.3.3. Why does $a - 1$ correspond to the number of prior successes, in other words, why is it not a , but $a - 1$?

Ex 8.1.28. BH.8.3.4.b. Given that the first patient is cured, what is the probability that the rest of the patients, i.e., the other $n - 1$, will also be cured?

Ex 8.1.29. Is this claim correct? Let T be the sum of two iid $\text{Unif}(0, 1)$ rvs. Then there exist a, b such that $T \sim \text{Beta}(a, b)$. (You don't need to derive the distribution of T .)

Ex 8.1.30. Show that $\beta(1, b) = 1/b$ by integrating the PDF of the beta distribution for $a = 1$. (Do not use the results of BH 8.5 for this exercise.)

Ex 8.1.31. Let $a, b > 1$. Show that the PDF of the beta distribution attains a maximum at $x = \frac{a-1}{a+b-2}$. Explicitly indicate where the assumption that $a, b > 1$ is used.

Ex 8.1.32. Explain in your own words:

1. What is a prior?
2. What is a conjugate prior?

Ex 8.1.33.

1. Look up on the web: what is the conjugate prior of the multinomial distribution? Give a name and a formula.
2. Explain why the Beta distribution is a special case of this distribution.

Ex 8.1.34. You make a test with n multiple choice questions and you give the correct answer to each question independently with probability p . The teacher's prior belief about p is reflected by a uniform distribution: $p \sim \text{Unif}(0, 1)$. Let X be the number of correct answers you give. What is the teacher's posterior distribution $p|X = k$? (You don't have to do a lot of math here; simply use a result from the book.)

Ex 8.1.35. You find a coin on the street. Initially, you are rather confident that this should be (approximately) a fair coin. This is reflected in your prior belief of the probability p of heads: $p \sim \text{Beta}(10, 10)$. Your friend is a bit more skeptical and assumes a uniform prior: $p \sim \text{Unif}(0, 1)$. You toss the coin 1000 times, and it comes up heads 900 times.

1. Determine your posterior distribution. (Again, use a result from the book)
2. Determine your friend's posterior distribution.
3. Compare the means of your posterior distribution and your friend's posterior distribution. Comment on the effect of the prior distribution if you have a lot of data.

Ex 8.1.36. We have an urn with 1000 coins. One of those is biased such that $P\{H\} = 99/100 = 1 - P\{T\}$, all others are fair. You select at random a coin, i.e., with probability $1/1000$ you select the biased one, and start throwing. You see 10 heads in row.

Ex 8.1.37. BH.8.3.5. write that X_j is an indicator of the j th throw beind made. Can this be the formal definition: $X_j = I_{N \geq j}$?

Ex 8.1.38. Use the pmf of the Beta-Binomial distribution to prove the following identity:

$$\sum_{k=0}^n \frac{\binom{n}{k} \binom{a+b-2}{a-1} (a+b-1)}{\binom{a+b+n-2}{a+k-1} (a+b+n-1)} = 1.$$

for all positive integers a, b, n .

Section 8.4

Ex 8.1.39. What is the SCV of Gamma(n, λ) distributed rv X ?

Ex 8.1.40. We have a machine that has temperature $x_0 e^{-\alpha t}$ after a time t . We switch it on when q jobs arrive. Job interarrival times are Exp(λ). Why does the temperature at the moment the q th job arrives have the distribution $x_0 \exp -\alpha Y$, with $Y \sim \text{Gamma}(q-1, \lambda)$?

Ex 8.1.41. Consider the chi-square distribution (with one degree of freedom) from BH.8.1.4.

Starting from the expression $f_Y(y) = \varphi(\sqrt{y}) y^{-1/2}$ in this example, show that this chi-square distribution is a special case of the Gamma distribution and specify the corresponding values of the parameters a and λ .

Ex 8.1.42. Is the sum of any two Gamma distributions again Gamma?

Ex 8.1.43. Prove by induction that $\Gamma(n) = (n-1)!$ if n is a positive integer.

Ex 8.1.44. Is the Poisson distribution the conjugate prior of the Gamma distribution?

Ex 8.1.45. Let $X \sim \text{Gamma}(4, 2)$ and $Y \sim \text{Gamma}(7, 2)$ be independent rvs. What is the distribution of $X + Y$? What is the distribution of $\frac{X}{X+Y}$?

Section 8.6

Read the definition of an order statistic. Skip the rest of BH.8.6.

Ex 8.1.46. If you can answer this question, then you basically know everything you need to know about order statistics for the purpose of this course.)

Let X_1, X_2, \dots, X_9 be a collection of random variables. Fill in the gaps (with just one word each time):

1. $X_{(1)}$ denotes the ... of X_1, X_2, \dots, X_9 .
2. $X_{(9)}$ denotes the ... of X_1, X_2, \dots, X_9 .
3. $X_{(5)}$ denotes the ... of X_1, X_2, \dots, X_9 .

8.2 BH EXERCISES: HINTS AND SOLUTIONS

Ex 8.2.1. BH.8.11. With convolution we know how to add and subtract independent rvs. Now we make a start with division. You'll see that this operator is not as simple as you always thought.

Before solving the problem, let's take a step back. You learned arithmetic at primary school. In all those problems, the numbers you had to add, subtract, etc. were supposed to be known precisely. At secondary school, you learned how to arithmetic with symbols. And now, at university, your next step is learn how to do arithmetic with rvs.

Here is an example to show you the relevance of this. In a paint factory at which a couple of my students did their master's thesis, the inventory level of dyes and other raw materials is often not known exactly. There are plenty of simple explanations for this. Raw materials are kept in big bags, and personnel uses shovels to take it out of the bags. Of course, occasionally, there is some spillage on the floor, and this extra 'demand' is not reported. The demand side is also not exact. A customer orders for example 500 kg of red paint. To make this, the operators follow a recipe, but dyes (in certain combinations) do not always give the same result. Therefore, the paint for each order is checked, and when it does not meet the quality level, the batch has to be adjusted by adding a bit more of certain dyes or solvents, or other chemical products.

When the planner has to make a decision on when to reorder a certain raw material, s/he divides the total amount of raw material by the average demand size. And this leads to occasional stock outs. When the stock level and the demands are treated as a rvs, such stock outs may be prevented, but this requires to be capable of determining the distribution of the something like Y/X .

Ex 8.2.2. BH.8.15. We'll use this exercise in a lecture to show how the normal distribution originates from astronomy (or dart throwing).

The notation is a bit clumsy for the angle coordinate. Write Θ for the rv and θ for its value.

Ex 8.2.3. BH. 8.18. Here we deal with division of rvs.

Ex 8.2.4. BH.8.23. We already analyzed how to handle addition, subtraction and division. It remains to deal with multiplication.

Ex 8.2.5. BH.8.31

Ex 8.2.6. BH.8.36.

Ex 8.2.7. BH.8.40. A nice question on the exam could be to take another prior, e.g., p uniform on $[1/3, 2/3]$. How would that affect the solution?

Ex 8.2.8. BH.8.52. The concepts discussed here are useful to better understand how to generate exponential random numbers.

Ex 8.2.9. BH.8.54. We tackle this also with simulation in an assignment.

I find it easier to consider $Y = pX$, rather than pX/q . Note that since $q = 1 - p \rightarrow 1$ as $p \rightarrow 0$, the factor $1/q$ is immaterial for the final result.

Read my solution too, as I develop some nice ideas in passing.

8.3 CHALLENGE: PING PONG BALLS IN A BELUGA

This challenge is a continuation of the simulation we did for the Beluga case, and we discuss some ways to check whether $V[N] \approx V[V]V[v]$ holds in general, and then we try to find a better approximation. We chopped up the challenge into many exercises, to help you organize the ideas.

Recall that earlier we have been a bit sloppy about the units, measuring the volumes of the airplane in m^3 and a ping pong ball in cm^3 , so actually N is in millions of ping pong balls. Note that using different units can easily lead to confusion; as a take-away, choose one unit.

One way to check the correctness of $V[N] \approx V[V]V[v]$ is to change the scale. In fact, memorize that changing scale is an easy way to check laws.

Ex 8.3.1. Suppose we instead measure the size of a ping pong ball in meters and the size of the airplane in hectometers. Explain that N is still in millions of ping pong balls. What happens to $V[N]$ and what happens to $V[V]V[v]$ (theoretically)?

Another way to check a statement is to consider some extreme cases.

Ex 8.3.2. Suppose that we would know the size of a ping pong ball very accurately, i.e. we consider the extreme case where $V[v] \rightarrow 0$. Explain that the approximation is not a good approximation in this limit.

Ex 8.3.3. Which of these two checks convinces you most that something is wrong with this approximation, and why?

We now turn to the task of trying to find a good approximation.

Ex 8.3.4. Assume that X and Y are independent. Show that

$$V[XY] = V[X]V[Y] + V[X]E[Y]^2 + E[X]^2V[Y].$$

Ex 8.3.5. Assume in addition that we know at least one of X and Y quite precisely. Argue that the following is then a good approximation:

$$V[XY] \approx V[X]E[Y]^2 + E[X]^2V[Y].$$

So far we have only considered the variance of a product, but we would like to know the variance of a ratio. For this we can use Taylor expansions to make accurate approximations.

Ex 8.3.6. Find the first order Taylor expansion of $\frac{1}{Z}$ around $a = E[Z]$. By taking the expectation and the variance of this expansion, show that

$$E\left[\frac{1}{Z}\right] \approx \frac{1}{E[Z]}, \quad V\left[\frac{1}{Z}\right] \approx \frac{V[Z]}{E[Z]^4}.$$

Ex 8.3.7. Combine all of the above to derive the following approximation for the variance of the ratio of two independent random variables X and Z :

$$V\left[\frac{X}{Z}\right] \approx \frac{V[X]}{E[Z]^2} + E[X]^2 \frac{V[Z]}{E[Z]^4}.$$

Ex 8.3.8. Check this approximation in the ways of the first two exercises.

After doing all this work, we would of course like to know how well this approximation does. When comparing the approximation to the sample standard deviation found in [7.1.22] for `num=500`, the result may be a bit disappointing. However, this is just because the sample standard deviation is also an estimate of the actual standard deviation of N , so by chance the result may be closer to $V[V]V[v]$ than to our new approximation.

In Chapter 10, you will learn something about the distribution of the sample variance. For now, just increase `num`. We know this decreases the variance of the sample mean and it also decreases the variance of the sample variance, so we get a more accurate estimate.

Ex 8.3.9. Use the result of the previous exercise to compute an approximation for $V[N] = V[V/v]$. Also use the code with a (much) higher value of `num`, to show that the approximation $V[N] \approx V[V]V[v]$ is likely to be worse, even in the setting of [7.2.15] where it was quite good.

The following two exercises are really optional, but I found them very neat and insightful.

Ex 8.3.10. Recall that for a non-negative random variable X with finite variance, we define the squared coefficient of variation as $SCV(X) = V[X]/E[X]^2$. Using the SCV, show that the approximations of [8.3.5] and [8.3.6] can be rewritten in the following neat way:

$$\begin{aligned} SCV(XY) &\approx SCV(X) + SCV(Y). \\ SCV(1/Z) &\approx SCV(Z). \end{aligned}$$

In BH.10, you will learn Jensen's inequality, which implies that $E\left[\frac{1}{Z}\right] \geq \frac{1}{E[Z]}$ for all positive random variables Z . In the following exercise, we reflect on this by finding a more accurate approximation based on the second order Taylor expansion.

Ex 8.3.11. Find the second order Taylor expansion of $\frac{1}{Z}$ around $a = E[Z]$. By taking the expectation, show that

$$E\left[\frac{1}{Z}\right] \approx \frac{1}{E[Z]} + \frac{2V[Z]}{E[Z]^3}.$$

Note that this is always at least $\frac{1}{E[Z]}$.

8.4 CHALLENGE: BENFORD'S LAW

In this exercise, we discuss Benford's law. Recall that the first step towards this law was taken in Lecture 5, Exercise 3. In this exercise, we showed that if X, Y are iid uniform on $[1, 10)$, that then the density of $Z = XY$ is given by

$$f_Z(z) = \frac{\log(\min\{10, z\}) - \log(\max\{1, z/10\})}{81} I_{1 \leq z \leq 100}.$$

Note that \log denotes the natural logarithm.

Benford's law is a statement about the distribution of the first digit of the product of sufficiently many variables that are iid uniform on $[1, 10)$. We first consider the first digit of the product of two such variables, i.e. the first digit of Z .

Ex 8.4.1. Let K be the first digit of Z . Show that the PMF of K is given by

$$P(K = k) = \frac{9k \log(k) - 9(k+1) \log(k+1) + 9 + 10 \log(10)}{81}$$

for $k \in \{1, 2, \dots, 9\}$.

Ex 8.4.2. Check that $\sum_{k=1}^9 P(K = k) = 1$. This can be done nicely by recognizing a *telescoping sum*: many terms cancel because they appear once with a minus and once with a plus.

Another way to derive the first digit of Z is to first divide Z by 10 if $Z \geq 10$. This yields a random variable W with support $[1, 10)$. Clearly, the division doesn't affect the first digit. The next exercise asks to derive the resulting density. This can be a bit tricky; you should check your answer by verifying that the distribution of the first digit of W matches the distribution of the first digit of Z .

Ex 8.4.3. Let $W = Z$ if $1 \leq Z < 10$ and $W = \frac{Z}{10}$ if $10 \leq Z < 100$. Derive the density of W .

We now turn to the product of more than two (independent) random variables. It would be very tedious to do this analytically, so we will instead use some code. However, to do this we have to approximate the continuous uniform variable by a discrete random variable. We use the discrete uniform distribution on $\{1 + 0.5 \cdot \frac{9}{s}, 1 + 1.5 \cdot \frac{9}{s}, 1 + 2.5 \cdot \frac{9}{s}, \dots, 1 + (s - 0.5) \cdot \frac{9}{s}\}$; in total this set has s elements. However, a product of two elements from this set may not again be an element of this set. To solve this, we identify all elements of the interval $(1 + k \cdot \frac{9}{s}, 1 + (k+1) \cdot \frac{9}{s})$ with $1 + (k+0.5) \cdot \frac{9}{s}$. We now use a loop to approximate the distribution of the product of $p+1$ random variables by looking at all possible values of the product of p random variables and one additional uniformly distributed random variable. Note that in the code, s is called `steps` and p is called `p_idx`.

Executing the code may take a while. If it takes more than 1 minute, you may decrease steps, but please do note that you did so.

Python Code

```
1 import math
```

```
2
```

```

3  steps = 900
4  products = 15
5  p_unif = [1.0/steps] * steps
6  p_mat = [p_unif]
7
8  for p_idx in range(1, products):
9      p_vec = [0] * steps
10     for s1 in range(steps):
11         for s2 in range(steps):
12             product = (1 + (s1 + 0.5)*9/steps) * (1 + (s2 + 0.5)*9/steps)
13             prod_probability = p_mat[p_idx - 1][s1] * 1/steps
14
15             if product > 10:
16                 product = product/10
17
18             prod_idx = math.floor((product-1)/9 * steps)
19             p_vec[prod_idx] += prod_probability
20
21     p_mat.append(p_vec)
22
23
24 p_digits = []
25 for p_idx in range(products):
26     vec = []
27     for digit in range(1, 10):
28         pd = sum(p_mat[p_idx][((digit-1)*steps//9):(digit*steps//9)])
29         vec.append(round(pd, 6))
30     p_digits.append(vec)
31
32 print(p_digits)

```

R Code

```

1  steps <- 900
2  products <- 15
3  p_unif <- rep(1/steps, steps)
4  p_mat <- matrix(0, nrow = steps, ncol = products)
5  p_mat[, 1] <- p_unif
6
7  for (p_idx in 2:products) {
8      p_vec <- rep(0, steps)
9      for (s1 in 1:steps) {
10         for (s2 in 1:steps) {
11             product <- (1 + (s1 - 0.5)*9/steps) * (1 + (s2 - 0.5)*9/steps)
12             prod_probability <- p_mat[s1, p_idx - 1] * 1/steps
13

```



```

14     if (product > 10) {
15         product <- product/10
16     }
17
18     prod_idx <- ceiling((product-1)/9 * steps)
19     p_vec[prod_idx] <- p_vec[prod_idx] + prod_probability
20 }
21 }
22 p_mat[, p_idx] = p_vec
23 }
24
25 p_digits <- matrix(0, nrow = 9, ncol = products)
26 for (p_idx in 1:products) {
27     for (digit in 1:9) {
28         pd = sum(p_mat[((digit-1)*(steps/9)+1):(digit*(steps/9)), p_idx])
29         p_digits[digit, p_idx] = round(pd, 6)
30     }
31 }
32 p_digits

```

Ex 8.4.4. Explain line P.13 (R.12) of the code.

Ex 8.4.5. Briefly comment on the results for $p = 2$ compared the exact result derived in the first exercise. Why is it important to make this comparison?

When looking at the results for larger p , it seems that the probabilities converge. The limit random variable B then satisfies the property that the first digit of B and the first digit of BU (where $U \sim \text{Unif}(1, 10)$) are identically distributed. Proving this is quite challenging (even for the challenge). In addition, we first need to know what the distribution of B is.

To guess the distribution of the first digit of B , we look at the results of our code and try some transformations to see if this yields familiar numbers. It turns out that the first digit M of B has the following distribution:

$$P(M = k) = \log_{10} \left(\frac{k+1}{k} \right),$$

for $k \in \{1, 2, \dots, 9\}$.

Ex 8.4.6. Briefly comment on these exact values of $P(M = k)$ compared to the values for $p = 15$ that result from the code. Give two reasons why the code results are not exact. Which reason do you think is the most important?

Ex 8.4.7. Check that $\sum_{k=1}^9 P(M = k) = 1$. You can again use a telescoping sum.

Besides the theoretical aspects covered in this challenge, Benford's law states that the first digit of numbers of naturally occurring sets that span several orders of magnitude, such as vote counts by county (or municipality), transaction sizes, etc., approximately follow this distribution. Initially this was just seen as an interesting curiosity of no practical value, but recently it has been used in fraud detection. If you're interested, you might check out this YouTube video by Numberphile: https://www.youtube.com/watch?v=XXjLR20K1kM&ab_channel=Numberphile

CHAPTER 9: QUESTIONS AND REMARKS

9.1 SIMPLE QUESTIONS

Section 9.1

Remark 9.1.1. Skip BH.9.1.6.

Ex 9.1.2. BH.9.1.7. We will also simulate this in an assignment.

I like this example as it shows how to make optimal decisions under uncertainty, but I have to admit that I don't understand the reasoning, or the use of conditional probability to solve this problem. Here is how I would solve the problem.

1. Explain that $X = (V - b)I_{b \geq \alpha V}$, with $\alpha = 2/3$, is the rv that models our payoff.
2. Why is this wrong: $E[X] = E[(V - b)I_{b \geq \alpha V}] = V E[I_{b \geq \alpha V}] - b E[I_{b \geq \alpha V}]$?
3. Compute $E[X]$, and provide a bound on α to ensure that $E[X] > 0$.

Ex 9.1.3. BH.9.1.7. Suppose $b < \alpha V$, for some $\alpha \in (0, 1)$ the bid is rejected. By the solution of BH.9.1.7, if $\alpha = 2/3$, you should not enter this game. But is there another α (smaller or larger than $2/3$) at which entering the game is interesting?

Ex 9.1.4. BH.9.1.8. Apply the same type of argumentation to find $E[X]$ when $X \sim \text{FS}(p)$.

Ex 9.1.5. BH.9.1.8. Use first step analysis to find $N_r := E[X]$ when $X \sim \text{NBin}(r, p)$.

Ex 9.1.6. BH.9.1.9. I reason slightly differently here. Write N_r for the number of throws required to reach r heads in row. Then I need N_{r-1} throws in expectation to reach the state in which there are $r - 1$ heads in row. Suppose now that we are in this state, i.e., there are $r - 1$ heads in row. Then, if I throw heads, with probability p , I reach the state with r heads in row, and I am done. However, if I throw tails, with probability q , I have to start all over again. Use this argument to derive the recursion $N_r = N_{r-1} + p \cdot 1 + q(1 + N_r)$. Solve this to obtain

$$N_r = \sum_{i=1}^r 1/p^i. \quad (9.1.1)$$

Ex 9.1.7. Compute the expected outcome of a die throw (with a 6-sided die), given that the outcome is even. Introduce proper notation for random variables and events.

Ex 9.1.8. BH.9.1.10. Let p_i , $0 \leq i \leq b$ be the probability to hit b before 0 .

1. Why is $p_0 = p_1/2$?

2. Why is $p_b = 1$?
3. Explain the recursion $p_i = p_{i-1}/2 + p_{i+1}/2$ for $0 < i < b$?
4. Show that point 3 implies that $p_i = \alpha i$ for $0 < i < b$ for any α we like.
5. Combine the fact that $p_b = 1$ with $p_i = \alpha i$ for all $0 < i < b$ to see that $\alpha = 1/b$.
6. Conclude that $p_0 = 1/2b$.

Section 9.2

Remark 9.1.9. About BH. 9.2.1. This definition is subtle, and it takes time to understand. Here is a slightly different explanation; perhaps it's useful for you.

Take some random variable X , say. Then, as in Chapter 7, we can be interested in $E[g(X)]$, i.e., the expectation of the rv $g(X)$.

When Y is continuous we can compute $E[Y | X = x]$ with the conditional CDF

$$E[Y | X = x] = \int f_{Y|X}(y|x) dy.$$

(For discrete rv., replace the integral by the PMF.) Observe that this is just a function of x ; define this function as $g(x) = \int f_{Y|X}(y|x) dy$. And now, as before, we consider the random variable $g(X)$, and we *call this rv the conditional expectation* of Y given X .

It is true that X plays some sort of double role—first we use it in the conditioning in the definition of the function g , and then we plug it into g again—and this is perhaps confusing. But I finally ‘got it’, when I understood that g can be interpreted as just some function of x . And then we compute $E[g(X)]$, and so on.

Ex 9.1.10. BH.9.2.2. Is $E[Y | I_A]$? a number or a rv?

Section 9.3

Remark 9.1.11. On BH.9.3.2 (Taking out what is known.) Perhaps it is easier to cristalize X into x . Then $g(x) = E[h(x)Y | X] = h(x)E[Y | X]$, because $h(x)$ is just a function. The rvs $E[H(X)Y | X]$ and $h(X)E[Y | X]$ are then both equal to $g(X)$.

Ex 9.1.12. On BH.9.3.9. Show that $\text{Cov}[Y - E[Y | X], E[Y | X]] = 0$.

Section 9.4

Ex 9.1.13. Consider a casino where, for any $a > 0$, it is possible to pay a euro and get a chance of $\frac{1}{5}$ on receiving $4a$ euro and a chance of $\frac{4}{5}$ of receiving nothing. Adam enters the casino with b euros, and bets half of his money on this gamble. Let X be the amount of money he has after the gamble. After that, he again bets half of the money he then has (i.e. half of X) on this gamble. Let Y be the amount of money he has after the second gamble.

1. Compute $E[X]$.
2. Compute $E[Y|X]$.
3. Compute $E[Y]$.

Explicitly mention the laws/rules you use.

Ex 9.1.14. Let $N \sim \text{Pois}(\lambda)$, and let $X|N \sim \text{Bin}(N, p)$, where $p \in (0, 1)$ and $\lambda > 0$ are known constants. Compute $E[X]$ using Adam's law. Check your answer using the chicken-egg story; with this story you can also obtain the distribution of X .

Ex 9.1.15. Correct? If A is an event and I_A is its indicator, then for all random variables X we have $E[X|A] = E[X|I_A]$.

Ex 9.1.16. Correct? If X and Y are independent, then $V[E[Y|X]] = 0$.

Ex 9.1.17. Let $X \sim \text{Exp}(\lambda)$, and let a be a constant.

1. Compute $E[X|X \geq a]$ using an integral and an indicator.
2. Explain the answer using a property of the exponential distribution.

Ex 9.1.18. A hat contains 9 fair coins and one coin that lands heads with probability 0.8. You pick a coin from the hat uniformly at random and toss it 10 times. Let A be the event that you pick a fair coin, and let X be the number of heads. Let B be the event that the first four tosses all show heads.

1. Compute $E[X|A]$.
2. Compute $E[X|A^c]$.
3. Compute $E[X]$.
4. Compute $P\{B\}$.
5. Compute $P\{A|B\}$.
6. Compute $E[X|B]$.
7. Compute $E[X|B^c]$. *Hint:* it is not necessary to compute $P\{A|B^c\}$.

Ex 9.1.19. Consider random variables $X, Y \in [0, 1]^2$ with joint PDF $f_{X,Y}(x, y) = 2I_{x \leq y}$. Determine $E[Y|X]$ and $E[X|Y]$.

Ex 9.1.20. Prove that $E[X|X \geq a] > E[X]$ for any a with $0 < P\{X \geq a\} < 1$.

Ex 9.1.21. Let $N \sim \text{Pois}(\lambda)$ and let $X|N \sim \text{Bin}(N, p)$, where $p \in (0, 1)$ and $\lambda > 0$ are known constants. Find $E[N|X]$.

Section 9.5

Ex 9.1.22. BH.9.5.1. Is $V[Y|X] = V[E[Y|X]]$?

Ex 9.1.23. Use Eve's law to show that $V[Y] \geq V[E[Y|X]]$.

Ex 9.1.24. Let $Z \sim \mathcal{N}(\mu, \sigma^2)$ and $Y = \sqrt{Z} + Z^2$. Find $V[Y|Z]$.

Ex 9.1.25. Correct? $V[Y] = V[Y|A] P\{A\} + V[Y|A^c] P\{A^c\}$ for any random variable Y and event A .

Ex 9.1.26. Let X, Y be random variables. Explain the difference between $V[Y|X]$ and $V[Y|X = x]$.

Ex 9.1.27. Show that $E[(Y - E[Y|X])^2|X] = E[Y^2|X] - (E[Y|X])^2$.

Ex 9.1.28. Let $X \sim \mathcal{N}(\mu, \sigma^2)$ and $W|X \sim \mathcal{N}(0, X^2)$. Find $V[W]$.

9.2 BH EXERCISES: HINTS AND SOLUTIONS

Ex 9.2.1. BH.9.1. It is easy to associate an indicator with the route chosen: let $R \in \{1, 2, 3\}$ be the rv corresponding to the route.

We tackle this problem also in an assignment.

Ex 9.2.2. BH.9.25. We tackle this problem also in an assignment with simulation. Check out https://en.wikipedia.org/wiki/Kelly_criterion you're interested.

Ex 9.2.3. BH.9.28.

Ex 9.2.4. BH.9.32. The results of this exercise are (or should be) used by nearly all software packages to control inventory levels of companies such as supermarkets and bol.com.

Ex 9.2.5. BH.9.37.

Bootstrapping is used in statistics to, for instance, construct confidence intervals. It is a much used and intuitive technique.

Extra exercise to help you recall some ideas of Ch 1. How many different bootstrap samples are possible?

I used some extra ideas to save some time. We say that the rvs $\{X_i\}$ are independent and distributed as the common rv X when $X_i \sim F_X$ where F_X is the CDF of the rv X . Then $E[X_i] = E[X]$, and so on. Next, I prefer to write $Y_j = X_j^*$, as this writes (and types) faster. Finally, it is easy to define $Y_j = \sum_{i=1}^n X_i I_{S_j=i}$, where $S_j \sim \text{DUnif}(\{1, \dots, n\})$ is the j th sample of the $\{X_i\}$.

Ex 9.2.6. BH.9.39. There are numerous examples of rvs with non-zero kurtosis, for instance, claim sizes of car accidents, the time patients spend in hospital beds, finance. This exercise helps to understand how a positive kurtosis may originate.

Ex 9.2.7. BH.9.50. We will also simulate this in an assignment.

Ex 9.2.8. BH.9.52

Ex 9.2.9. BH.9.55. Suppose first you draw just one number per day, what is then the recursion? Then suppose you draw 2 numbers per day.

An interesting variation is to find a recursion for the number of *draws* instead of *days* are needed until all numbers have been seen.

Ex 9.2.10. BH.9.56.

Ex 9.2.11. BH.9.57

Ex 9.2.12. BH.9.58. In part c. the prior is the uniform distribution. What would happen if you would take the prior of part b, i.e., a out of j wins?

9.3 CHALLENGE: BETTING

Consider the setting of BH.9.25, which you also studied in the coding section. We use the notation from that exercise. In this exercise we will discuss how to set f , the betting fraction. In particular, we will discuss the *Kelly criterion*, which states that the betting fraction should be $f = 2p - 1$ if $p > \frac{1}{2}$ is the winning probability.

We discuss its relationship to expected utility theory, which you will also study in Introduction to Mathematical Economics. Expected utility theory states that bets should be chosen to maximize expected utility. So we solve $\max_{0 \leq f \leq 1} E[U(X_{n+1})|X_n]$.

Ex 9.3.1. Show that solving the maximization problem for the utility function $U(x) = \log(x)$ yields the betting fractions from the Kelly criterion, $f = 2p - 1$ if $p > \frac{1}{2}$ and $f = 0$ if $p \leq \frac{1}{2}$.

Other people may have a different utility function, which yields a different betting fraction.

Ex 9.3.2. Calculate the utility maximizing betting fraction f if $U(x) = \sqrt{x}$.

Note that for both of these utility functions, the betting fraction f does not depend on the wealth X_n before the gamble, but in general f does depend on X_n .

Now that we have two different betting fractions, we compare them. For that, we first need the following result:

Ex 9.3.3. Assume that f does not depend on X_n . Let $x_0 = 1$. Show that there exist constants a, b such that $\log(X_n) = aW + b$ and $W \sim \text{Bin}(n, p)$, and determine a and b in terms of f .

Theorem 10.3.6. states that (for sufficiently large n) we can approximate a random variable with the binomial distribution $W \sim \text{Bin}(n, p)$ by a random variable with the normal distribution $\text{Norm}(np, np(1-p))$. While you will only learn about the proof of this next week, we are already going to use this approximation here.

Ex 9.3.4. Two people (Carl and Daria) participate in n rounds of this betting game. Their games are independent. Carl's initial wealth is $x_0 = 1$ and Daria's initial wealth is $y_0 = 1$. We denote Carl's wealth after n rounds by X_n and Daria's wealth after n rounds by Y_n . Carl chooses f according to the Kelly criterion, i.e. $f = 2p - 1$. Daria chooses f to be the utility maximizing betting fraction for $U(x) = \sqrt{x}$. Use the previously mentioned normal approximation to derive an approximation for the difference $\log(Y_n) - \log(X_n)$.

Kelly's criterion does not mention utility functions, it just recommends to set $f = 2p - 1$ regardless of one's utility function. The next exercise is meant to give some insight why.

Ex 9.3.5. Use `pnorm` in R, or `norm.cdf` in Python, to approximate $P(X_n > Y_n)$ for some chosen values for n and p . What do you think that happens if $n \rightarrow \infty$ for a fixed p ? Explain why this is an argument to use the Kelly criterion regardless of one's utility function. Also, explain why maximizing utility suggests a different f in spite of this result.

CHAPTER 10: QUESTIONS AND REMARKS

10.1 SIMPLE QUESTIONS

Section 10.1

Ex 10.1.1. On BH.10.1.1: here is perhaps simpler proof of the Cauchy-Schwarz inequality. Define $f(t) = E[(Y - tX)^2]$.

1. Explain that $f(t) \geq 0$.
2. Write $f(t) = E[(Y - tX)^2]$ as a polynomial of the second degree, i.e., in the form $f(t) = at^2 + bt + c$ (Hint, see the proof of BH.10.1.1).
3. Since $f(t) \geq 0$, how many (real) roots can it have at most?
4. What are the implications of this for the discriminant $D = b^2 - 4ac$?
5. Show that the Cauchy-Schwarz inequality directly follows from this restriction on D .

Remark 10.1.2. I find it easier to remember the Cauchy-Schwarz inequality in the form $(E[XY])^2 \leq E[X^2] E[Y^2]$; like this there are squares on both sides.

Ex 10.1.3. On BH.10.1.3. How do they get from $P\{X > 0\}$ to the inequality for $P\{X = 0\}$? (Provide the details.)

Ex 10.1.4. On BH.10.1.3. Do the algebra to show that $P\{X = 0\} = 1/(\mu + 1)$.

Ex 10.1.5. On BH.10.1.3. Explain that we actually use Markov's inequality.

Ex 10.1.6. On BH.10.1.3. What is the probability of two people with birthdays 2 days apart?

Remark 10.1.7. I often forget the direction in Jensen's inequality. To check, the following reasoning works for me: I know that $V[X] \geq 0$, but $V[X] = E[X^2] - (E[X])^2 = E[g(X)] - g(E[X])$ with $g(x) = x^2$. Then, from the graph of the parabola, i.e., the graph of g , I know that g is convex.

Ex 10.1.8. In Jensen's inequality, when does equality hold? Can you explain (in terms of convexity and concavity) why equality holds for only this type of functions?

Remark 10.1.9. Skip BH.10.1.7, 10.1.8, 10.1.9

Ex 10.1.10. When X is a non-negative rv, prove the simplest form of Markov's inequality: $P\{X \geq a\} \leq E[X]/a$ for $a \geq 0$. Then show that BH.10.1.10 follows from this.

Ex 10.1.11. Which of the following are equivalent to Chebyshev's inequality? Show why or why not.

1. $P(|X - E[X]| \geq a) \leq \frac{V[X]}{a^2}$ for all $a > 0$
2. $P(|X - E[X]| < a) > \frac{V[X]}{a^2}$ for all $a > 0$
3. $P(|X - E[X]| < a) \geq 1 - \frac{V[X]}{a^2}$ for all $a > 0$
4. $P(|X - E[X]| \geq c\sigma) \leq \frac{1}{c^2}$ for all $c > 0$ and $\sigma^2 = V[X]$.

Ex 10.1.12. On BH.10.1.11. Why is Chebyshev's inequality of no use if we try to plug in values for $0 < a \leq \sqrt{V[X]}$?

Ex 10.1.13. BH.10.1.13 shows that Chernoff's inequality is a very strict bound. Is Chernoff's inequality always the tightest bound (out of the ones you know)? What about the case where X is defined as follows

$$P\{X = 0\} = \frac{3}{4}, \quad P\{X = 2\} = \frac{1}{4}.$$

Section 10.2

Ex 10.1.14. Is the following statement equivalent to the strong or the weak law of large numbers? Fix $\epsilon > 0$. For all $\delta > 0$, there is an n so large that $P\{|\bar{X}_n - \mu| > \epsilon\} \leq \delta$.

Ex 10.1.15. Which of the two following statements correctly represents the strong law: $\lim_{n \rightarrow \infty} P\{|X_n - \mu| > \epsilon\} = 1$ for all $\epsilon > 0$, or $P\{\lim_{n \rightarrow \infty} |X_n - \mu| = 0\} = 1$.

Ex 10.1.16. On BH.10.2.5. Where have we applied this idea earlier?

Section 10.3

Ex 10.1.17. On BH.10.3.1. I prefer to write $\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$. Why could that be?

Ex 10.1.18. On BH.10.3.7. $E[\log Y_n] = \log 100 - 0.081n$. Explain the 0.081. Think also about the paradoxical outcome. (Once again, probability theory is hard.)

Ex 10.1.19. On BH.10.3.7. The stock rises α % and decreases β %. Find a relation between α, β such that $\lim Y_n \geq 1$.

Ex 10.1.20. On BH.10.3.7. Note that $E[\log Y_n] \sim -0.081n$, i.e., has negative drift, but $\log E[Y_n] \sim n \log 1.1$. Check that this is not in conflict with Jensen's inequality.

Section 10.4

Ex 10.1.21. On BH.10.4.3. Why is $n\bar{Z}_n^2 \sim \chi_1^2$?

Ex 10.1.22. On BH.10.4.3. Show that $\sum_j^n (Z_j - \bar{Z}_n)^2$ and \bar{Z}_n^2 are independent.

Ex 10.1.23. A fair coin is tossed 100 times. We are interested in the probability that the number of heads that turn up is at most 40. What is the tightest upper bound on this probability that we can find by using Chebyshev's inequality? Hint: use a symmetry argument.

Ex 10.1.24. Here is inequality from which all inequalities in BH 10.1.3 immediately follow. It's worth memorizing. Take any rv X and a function f that is non-negative and non-decreasing.

1. Why is this true for any a : $f(a)I_{X \geq a} \leq f(X)I_{X \geq a} \leq f(X)$?
2. Take expectations in the inequality of the previous step and use the fundamental bridge to show that $P\{X \geq a\} \leq E[f(X)]/f(a)$.
3. What part of the proof goes wrong if f can also be negative?
4. Show that Markov's inequality follows by taking $Y = |X|$ and $f(x) = x$. Why don't we take $f(x) = |x|$?
5. Show that Chebyshev's inequality follows by taking $Y = |X - \mu|$ and $f(x) = x^2$.
6. Show that Chernoff's inequality follows by taking $f(x) = e^x$.

Ex 10.1.25. Let the set of r.v.s $\{X_k, k \geq 1\}$ be the outcomes of throws of a biased coin. We take $X_j = 1$ if the outcome is heads, and $X_j = 0$ if tails. Suppose $E[X_k] = \mu$ and $V[X_k] = \sigma^2$. Let $Y_j = \sum_{i=n_{j-1}+1}^{(n+1)j} X_i/n$, i.e., Y_j is the sample mean of the j batch of throws. Since $\{Y_j, j \geq 1\}$ are iid, take Y as the common r.v., i.e., $Y_j \sim Y$. What is a (frequentist) explanation of the statement $P\{|Y - \mu| > \epsilon\} \leq \sigma^2/n\epsilon$?

Ex 10.1.26. Interpret the WLLN in terms of the previous exercise.

Ex 10.1.27. In the setting of [10.1.25], the probability of a sequence of outcomes like this: H, T, H, T, H, T, \dots , i.e., a sequence in which the heads and tails alternate, has probability zero. However, $\sum_{i=1}^n I_{X_i=H}/n \rightarrow 1/2$. So, we have a sequence that occurs with probability zero, but still the average along the sequence has the proper limit. Doesn't this violate the SLLN?

10.2 BH EXERCISES: HINTS AND SOLUTIONS

Ex 10.2.1. BH.10.2**Ex 10.2.2.** BH.10.3 This is just a funny exercise, but I wonder whether it has a practical value.**Ex 10.2.3.** BH.10.6**Ex 10.2.4.** B.10.9**Ex 10.2.5.** BH.10.23.**Ex 10.2.6.** BH.10.26.**Ex 10.2.7.** BH.10.28. Note that standardized version of a rv X is $Y = (X - \mu)/\sigma$ where $E[X] = \mu$ and $V[X] = \sigma$.**Ex 10.2.8.** BH.10.30. The problem demonstrates a simple investment strategy. If you plan to work as a quant in finance or as an actuary, or if you play poker, or some similar game, such strategies should interest you naturally.**Ex 10.2.9.** BH.10.36.**Ex 10.2.10.** BH.10.39.

10.3 CHALLENGE: RECORDS

In BH.7.48 you looked at the number of records in high jumping. Let X_j be how high the j th jumper jumped. As in that exercise, we assume that X_1, X_2, \dots are iid. with a continuous distribution and say that the j th jumper sets a record if X_j is larger than X_i for all $1 \leq i \leq j-1$. Let X_i^* denote the i th record, i.e., the height of highest jump for the first i jumps. We write f_X and F_X for the PDF and CDF of the iid jumping heights, and X for a random variable with density f_X . Finally, we write G_X for the survivor function of X , i.e. $G_X(x) = 1 - F_X(x)$.

It is not necessary to know the solution of BH.7.48 to do the challenge. In the challenge, we instead look at the distribution of the i th record, the expectation of the i th record and the expected improvement of the record: $E[X_{n+1}^* - X_n^*]$.

Ex 10.3.1. Let $f_{X_{i+1}^*, X_i^*}$ be the joint PDF of the $(i+1)$ th and the i th record. Prove that

$$f_{X_{i+1}^*, X_i^*}(u, v) = \frac{f_X(u)}{G_X(v)} f_{X_i^*}(v) I_{u > v}.$$

Note that $X_1^* = X_1$. We now derive the density of X_2^* , and then proceed with the general case. These are challenging problems, be sure to check out the hints if you are stuck.

Ex 10.3.2. Prove that

$$f_{X_2^*}(u) = -f_X(u) \log(G_X(u)).$$

Ex 10.3.3. Prove that

$$f_{X_n^*}(u) = f_X(u) \cdot \frac{(-\log(G_X(u)))^{n-1}}{\Gamma(n)}.$$

In general case, it is hard to compute the integral of $uf_{X_n^*}(u)$ which is required to compute $E[X_n^*]$. For the exponential distribution however, it is still possible to find the result analytically.

Ex 10.3.4. Assume that $X \sim \text{Exp}(\lambda)$. Determine $E[X_n^*]$, and hence the expected improvement of the record: $E[X_{n+1}^* - X_n^*]$, using the PDF from the previous exercise.

When X is exponentially distributed, we can also use directly use the properties of the exponential distribution to determine $E[X_{n+1}^* - X_n^*]$.

Ex 10.3.5. In Exercise 6 of Assignment 5 you found an expression for $E[X|X \geq a]$ if $X \sim \text{Exp}(\lambda)$. Use this expression and Adam's law to determine $E[X_{n+1}^* - X_n^*]$ if $X \sim \text{Exp}(\lambda)$.

Ex 10.3.6. Is the exponential distribution a realistic model for record improvements? Why (not)? If not, why is it still good to look at this case? Explain briefly.

We now consider this model for other distributions as well. Although for many distributions finding analytical results is very difficult or impossible, it can still be interesting to make a plot for the record improvements for other distributions.

Ex 10.3.7. Assume that X has PDF $f_X(x) = 2xe^{-x^2}$ for $x > 0$, and $f_X(x) = 0$ otherwise. Plot $E[X_{n+1}^* - X_n^*]$ as a function of n . You may use code for computing the survival function and the expectation, although it is possible (but not recommended) to do it analytically.

OLD EXAM QUESTIONS

Ex 11.0.1. Inspired by BH.7.3.6, compute $V[M]$ and $V[L]$ by first computing $E[M^2]$ and $E[L^2]$.

Ex 11.0.2. Let (X, Y) follow a Bivariate Normal distribution, with X and Y marginally following $\mathcal{N}(\mu, \sigma^2)$ and with correlation ρ between X and Y .

1. Use the definition of a Multivariate Normal Distribution to show that $(X + Y, X - Y)$ is also Bivariate Normal.
2. Find the marginal distributions of $X + Y$ and $X - Y$.
3. Compute $\text{Cov}[X + Y, X - Y]$. Then, write down the expression for the joint PDF of $(X + Y, X - Y)$.

Ex 11.0.3. This is about the simplest model for an insurance company that I can think of. We start with an initial capital $I_0 = 2$. The company receives claims and contributions every period, a week say. In the i th period, we receive a contribution X_i uniform on the set $\{1, 2, \dots, 10\}$ and a claim C_i uniform on $\{0, 1, \dots, 8\}$.

1. What is the meaning of $I_1 = I_0 + X_1 - C_1$?
2. What is the meaning of $I_2 = I_1 + X_2 - C_2$?
3. What is the interpretation of $I'_1 = \max\{I_0 - C_1, 0\} + X_1$?
4. What is the interpretation of $I'_2 = \max\{I'_1 - C_2, 0\} + X_2$?
5. What is the interpretation of $\bar{I}_n = \min\{I_i : 0 \leq i \leq n\}$?
6. What is $P\{I_1 < 0\}$?
7. What is $P\{I'_1 < 0\}$?
8. What is $P\{I_2 < 0\}$?
9. What is $P\{I'_2 < 0\}$?
10. Provide an interpretation in terms of the inventory of rice, say, at a supermarket for I_1 and I'_1 .

Ex 11.0.4. Take $X \sim \text{Unif}\{-2, -1, 1, 2\}$ and $Y = X^2$. What is the correlation coefficient of X and Y ? If we would consider another distribution for X , would that change the correlation?

Ex 11.0.5. We have a machine that consists of two components. The machine works as long as both components have not failed (in other words, the machine fails when one of the two components fails). Let X_i be the lifetime of component i .

1. What is the interpretation of $\min\{X_1, X_2\}$?
2. If X_1, X_2 iid $\sim \text{Exp}(10)$ (in hours), what is the probability that the machine is still ‘up’ (i.e., not failed) at time T ?
3. Use the previous result to determine the distribution of $\min\{X_1, X_2\}$.
4. What is the expected time until the machine fails?

Ex 11.0.6. We have two rvs X and Y with the joint PDF $f_{X,Y}(x,y) = \frac{6}{7}(x+y)^2$ for $x,y \in (0,1)$ and 0 else. Also we consider the two rvs U and V with the joint PDF $f_{U,V}(u,v) = 2$ for $u,v \in [0,1], u+v \leq 1$ and 0 else.

1. Compute $P\{X+Y > 1\}$.
2. Compute $\text{Cov}[U, V]$.

(Hint: first draw the area over which you want to integrate, if this does not help check out the discussion board post on exercise 7.13a from the first Tutorial)

Ex 11.0.7. Let $U = X + Y$ and $V = X - Y$ where $X, Y \sim U[0,1]$ and independent. Show that

$$f_{U,V}(u,v) = \frac{1}{2} I_{|v| \leq u \leq 2-|v|}.$$

Ex 11.0.8. Let X and Y have PDF $f_{X,Y}$. Take $g(x,y) = (\min\{x,y\}, \max\{x,y\})$. Why is

$$f_{U,V}(u,v) = f_{X,Y}(u,v) + f_{X,Y}(v,u)?$$

Simplify to $f_{U,V}(u,v) = 2f(u)f(v)$ for the case X, Y iid with common PDF f .

Ex 11.0.9. Let X, Y be continuous rvs with CDF $F_{X,Y}(x,y) = (x-1)^2(y-2)/8$ for $x \in (1,3)$.

- a. Explain that $y \in (2,4)$ for F to be a proper CDF.
- b. What is $F(3,7)$?
- c. Determine the PDF.
- d. Compute $P\{2 < X < 3\}$
- e. Compute $P\{2 < X < 3, 2 < Y < 3\}$.
- f. Compute $P\{Y < 2X\}$.
- g. Compute $P\{Y \leq 2X\}$.

h. Compute $P\{Y < 2X, Y + 2X > 6\}$.

Ex 11.0.10. Consider the general case where we are given the relationship $U = V^4$ between the random variables U and V for $V \in (-3, 2)$. Explain why we cannot simply invoke the change of variables theorem.

Now imagine V following a uniform distribution on the given interval. Consider the given transformation on the intervals $(-3, 0)$ and $(0, 2)$ separately. Explain why this allows you to employ the change of variables theorem and find the distribution of U on these intervals. Finally combine these results (using indicator functions) and state the PDF of U (remember to adjust the domain for the indicator functions according to the transformations).

Ex 11.0.11. Let $U \sim \text{Unif}(0, \pi)$. Use BH.8.1.9 to show that $X = \tan(U)$ has the Cauchy distribution. Compare this exercise to BH.8.1.5.

YOU WALK INTO a bar and you find two people, Amy and Bob, playing a game of darts. Their game consists of several rounds, called *legs*, and the first person to win 7 legs wins the match. You have never seen Amy or Bob play before, so you don't know their strength. Denoting by p the probability that Amy wins a leg, your prior belief can be modeled by a uniform distribution: $p \sim \text{Unif}(0, 1)$. (Note: we assume that p remains constant during the entire match; even though your *beliefs* about p might change.)

Denoting by A a leg won by Amy and by B a leg won by Bob, the result of the first 10 legs is: AAABBAABAB. Your friend Charles is very confident that Amy will win the match and he offers you a bet: if Amy wins the match, you must pay €10 to Charles; if Bob wins the match, he must pay you €300. You are tempted to take the bet, but you want to do some calculations first.

Ex 11.0.12. Is the order in which Amy and Bob won their respective legs relevant for your posterior probability that Bob will win the match?

Ex 11.0.13. Let A_n denote the number of legs that Amy won out of a total of n legs. Express the result of the first 10 legs in terms of A_n

Ex 11.0.14. What is the distribution of $A_n|p$ (i.e., the distribution of A_n given the value of p)?

Ex 11.0.15. Find the posterior distribution of p after observing $A_n = k$.

Ex 11.0.16. According to your posterior belief about p , what is the probability that Bob wins the match? Express your answer in terms of beta functions. (Hint: Use the law of total probability.)

Ex 11.0.17. Using the expression

$$\beta(a, b) = \frac{(a-1)!(b-1)!}{(a+b-1)!} \quad (11.0.1)$$

for every positive integers a, b , compute the probability from the previous question as a number.

Ex 11.0.18. Assuming that you want to maximize your expected outcome, should you take the bet?

Ex 11.0.19. On BH.9.5.4. With Z and W as given in the example, show that $E[Z|W] = \rho$ and $V[Z|W] = 1 - \rho^2$.

Ex 11.0.20. Prove that

$$E[(Y - E[Y|X] - h(X))^2] = E[(Y - E[Y|X])^2] + E[(h(X))^2] \quad (11.0.2)$$

for all random variables X, Y and all functions h . Explain why this result implies that $E[Y|X]$ is the best predictor of Y based on X .

IN THE NEXT couple of problems we derive Eve's law in a slightly different way than BH.

Define $\hat{X} = E[X|Y]$ as an *estimator* of X and $\tilde{X} = X - \hat{X}$ as the estimation error.

Ex 11.0.21. Show that $E[\tilde{X}|Y] = 0$.

Ex 11.0.22. Prove that $E[\tilde{X}] = 0$. What does it mean that $E[\tilde{X}] = 0$?

Ex 11.0.23. Prove that $E[\tilde{X}\hat{X}] = 0$.

Ex 11.0.24. Show that $\text{Cov}[\hat{X}, \tilde{X}] = 0$. Conclude that

$$V[X] = V[\hat{X} + \tilde{X}] = V[\hat{X}] + V[\tilde{X}]. \quad (11.0.3)$$

Ex 11.0.25. Prove that $V[\tilde{X}] = E[V[X|Y]]$. Conclude Eve's law.

CHEBYSHEV'S INEQUALITY IS useful in proving notions of convergence in probability, which you will see repeatedly in later courses. We say X_n converges in probability to the random variable Z if

$$\lim_{n \rightarrow \infty} P(|X_n - Z| \geq \varepsilon) = 0 \quad \forall \varepsilon > 0$$

Note that in the above definition setting $Z = a$ for some constant a would still be valid, as technically the constant a is a random variable.

Ex 11.0.26. Let Y_n denote the number of heads obtained from throwing a fair coin n times. Then $\frac{Y_n}{n}$ clearly is the proportion of heads in the sample. Find the expectation of this proportion, and show that it converges in probability to its mean. This is denoted as $\frac{Y_n}{n} \xrightarrow{P} E\left[\frac{Y_n}{n}\right]$ and is known as the weak law of large numbers.

Ex 11.0.27. Where would this proof break down if we try to apply it to e.g. the Cauchy distribution?

HINTS

h.7.1.16. Recall that $F \in [0, 1]$.

h.7.1.28. Realize that $E[ML] = E[XY]$.

h.7.2.6. For 2D LOTUS, take $g(i, j) = \min\{i, j\}$.

h.7.2.7. $X \sim \text{Geo}(1 - q) \implies E[X] = q/(1 - q)$. Now use that $L \sim \text{Geo}(1 - q^2)$.

h.7.2.8. Use that $1 - q^2 = (1 - q)(1 + q)$.

h.7.2.10. Use 2D LOTUS on $g(x, y) = I_{\max\{x, y\}=k}$.

h.7.2.12.

h.7.2.16. Marginalize out L .

h.7.2.17. Marginalize out M .

h.7.2.19. Use partial integration.

h.7.2.20. Using independence and the specific property of the r.v. L that $\{L > x\} \iff \{X > x, Y > x\}$:

h.7.2.28. If you recall the Poisson distribution, you know that $e^\lambda = \sum_{i=0}^{\infty} \lambda^i / i!$. In fact, this is precisely Taylor's expansion of e^λ .

h.7.2.29. Use [7.2.7].

h.7.2.31. The idea is not difficult, but the technical details require attention, in particular the limits in the integrations.

h.7.2.33. Realize that $P\{L = i, M - L = j\} = P\{L = i, M = i + j\}$. Now fill in the formula of [7.2.15].

h.7.2.36. Reread BH.7.2.2. to realize that $E[M - L] = E[|X - Y|]$. Relate the latter expectation to the expression in the problem.

h.7.2.37.

h.7.2.39. Observe that Z is independent from X and Y , hence from M and L

h.7.4.1. Check BH 7.2.2. Bigger hint: Let A the arrival time of Alice, and B the time of Bob. Then we want to compute $P\{|A - B| \leq 1/4\}$. (15 minutes is $1/4$ hour.) Why is $f_{A,B}(a,b) = I_{a \in [0,1]} I_{b \in [0,1]}$? Now apply 2D-LOTUS to the function $g(a,b) = I_{|a-b| \leq 1/4}$.

h.7.4.2. a. $P\{X = i, Y = j, N = n\} = P\{X = i, Y = j\} I_{i+j=n}$.
c. $P\{X = i | N = n\} = 1/(n+1)$. Why is this uniform?

h.7.4.4. a. First find $f_{Y|X}$ and $f_{Z|X}$. Then, given X , Z and Y are iid. Hence $f_{X,Y,Z} = f_{Y,Z|X} f_X$. Use independence to split $f_{Y,Z|X}$ into a product. b. Suppose that Y is really big. Since Y is dependent on X , X must be dependent on Y . But Z is in turn dependent on X . What are the consequences?

h.7.4.5. Section 7.2

h.7.4.6. Make a drawing.

h.7.4.7. Check BH.7.1.24 and BH.7.1.25 First draw the area over which we have to integrate. Then use an indicator function over which to integrate. What is the joint PDF $f_{Y_1, Y-2}$?

h.7.4.8. Section 7.2

h.7.4.10. Use the hint of the book and independence to see that $E[S_n^2 T_n^2] = E[S_n^2] E[T_n^2]$. Then try to simplify.

b. It is immediate that $E[S_n] = 0$. Hence, focus on $E[S_n T_n]$. Expand the sums of $E[S_n T_n]$, and consider the individual terms $E[X_i Y_j]$. When $i \neq j$, are X_i and Y_j independent? What if $i = j$?

c. It is clear that $R_n^2 = S_n^2 + T_n^2$. Now use linearity to split $E[R_n^2]$. Finally, realize that $E[S_n] = 0$, hence $E[S_n^2] = V[S_n]$. But then we can use the formula of the variance of a sum to split it up into a sum of variances plus covariances.

h.7.4.11. a. Expand the brackets in the expression for the sample variance r to see that

$$r = 1/n \sum_i x_i y_i - \bar{x} \bar{y}.$$

Next, we choose with probability $1/n$ one the points (x_i, y_i) . Under this probability, $E[XY] = 1/n \sum_i x_i y_i$, $E[X] = \bar{x}$, $E[Y] = \bar{y}$. So, how do $\text{Cov}[X, Y]$ and r relate?

b. Expand the brackets and use iid and linearity properties to show that the expected area spanned by two random points (X, Y) and (\tilde{X}, \tilde{Y}) satisfies

$$E[(X - \tilde{X})(Y - \tilde{Y})] = 2 \text{Cov}[X, Y].$$

h.7.4.12. a. Use that expectation is linear.

b. Read the entire exercise in its entirety before trying to solve it. In this case trying to solve c. seems simpler because of the extra iid assumption. You might want to use this to formulate some simple guesses.

Thus first part c. It is given that the X_i and Y_j are iid. Then, if I could improve the estimator $\hat{\theta}$ by splitting the measurements into two sets X_i and Y_j , then I would certainly do that. And not only I would do that; anybody in his right mind would do that. But, I never heard of this idea, and I am sure you have neither, so this must be impossible (because if it would, people would have been using this trick for ages.) Hence, we can place this in the context of the maxim: ‘we cannot obtain information for free’. For this case, this must imply that splitting iid measurements into smaller sets cannot help with improving the estimator. What does this idea imply for the weights?

Part b, continued. I always try to solve the problem myself without a hint. This lead to the following considerations, which gave me quite a bit of extra understanding beyond the problem itself. As a next piece of advice, before doing hard work, I prefer to look at some corner cases to acquire some intuitive understanding. I also use the rvs of Part c.

Suppose that $v_2 := V[Y_j] = 0$, but $v_1 := V[X_i] > 0$. (For instance, Y_j is the j th measurement of a perfect machine and X_j of an imperfect machine.) Then we know that the set $\{Y_j\}$ forms a set of perfect measurements. But then I am not interested in the $\{X_i\}$ measurements anymore; why should I as I have the perfect measurements $\{Y_j\}$ at my disposal. So, then I put $w_1 = 0$, because I don’t want the $\{X_i\}$ measurements to pollute my estimator. In other words, the final result should be such that $v_2 = 0 \implies w_1 = 0$, and vice versa.

More generally, I learned from this corner case that I want this for the final result: when $v_2 < v_1 \implies w_1 < w_2$, and vice versa.

How would you choose the weights such that this requirement is satisfied, but also the condition imposed by Part c.?

h.7.4.13. b. The people in the sample of size n with an A is $X_1 + X_2$. But this is the same as $n - X_3$. Hence, what is $P\{X_3 = n - i\}$?

c. I found this a hard problem. Here is my hint based on recursion. Let S_n be the number of A s in n individuals. We want to know $f_n(i) = P\{S_n = i\}$. A simple recursive idea, i.e., one-step analysis by conditioning on the phenotype of the n th person, gives that

$$f_n(i) = f_{n-1}(i-2)p^2 + f_{n-1}(i-1)2pq + f_{n-1}(i)q^2,$$

with $q = 1 - p$ as always. Now I was a bit stuck, but just to try to see whether I could see some structure, I tried a simpler case, namely, a recursion for the binomial distribution. Derive this, and then use this to solve the problem.

d. It is easiest to work with $f(p) = \log P\{X_1 = k, X_2 = l, X_3 = m\}$, where $P\{X_1 = k, X_2 = l, X_3 = m\}$ follows from a., and then differentiate with respect to p .

e. Follow the same scheme as for d.

h.7.4.14. The challenge for you is to try to understand the mathematics behind these concepts. Read the exercise a number of times. I found it quite difficult to capture the

concepts in formulas. (I solved it once. After two weeks, I tried to solve it again, and found it just as hard as the first time.) Once you have the model, the technical part itself is simple.

h.7.5.1. In this exercise we want to prove that N is Poisson distributed. So you cannot assume this in your solution.

h.7.5.3. Use the relation of the previous exercise to show that

$$P(N = n + 1) = \frac{\lambda}{1 + n} P(N = n). \quad (12.0.1)$$

Bigger hint: Fill in $y = 0$ in the LHS and RHS of (7.5.1); call this expression 1. Then fill in $y = 1$ to obtain a second expression. Divide these two expressions and note that $P\{X = x\}$ cancels. Finally, define

$$\lambda = \frac{P\{Y = 1\}}{(1 - p)P\{Y = 0\}}. \quad (12.0.2)$$

h.7.8.2. First, compute the value of $E[v(d, x)]$ as a function of x . Then find the optimal value of x .

h.7.8.5. You don't have to compute x^* for the case where $\rho = 1$; this is not easy!

h.8.1.8.

h.8.1.9.

h.8.1.10.

h.8.1.19. Let X, Y be iid standard normal. Since the square of a standard normal r.v. is chi-square distributed, we can write S as $S = X^2 + Y^2$ (here we use BH.8.1.4).

h.8.1.21.

h.8.1.22.

h.8.1.23.

h.8.1.24.

h.8.1.25.

h.8.1.26.

h.8.1.27.

h.8.1.28.

h.8.1.36.

h.8.1.37.

h.8.2.1. Start with the case $v = 0$. Use the proof of BH.8.1.1. Reason carefully; corner cases as simple to miss.

Then, make a graph of the two branches of the hyperbola's $1/t$, one branch for $t > 0$, the other for $t < 0$. Then draw a horizontal line to indicate the level $V = v$; this shows with part(s) of the hyperbola's lie below v . Then compute the probability for each branch. This will give the answer of the book immediately.

h.8.2.2. a. See BH.8.1.9.

b. If (X, Y) are uniform on the disk, then the function $g(x, y)$ must be constant on this disk. Use an indicator to ensure that $X^2 + Y^2 \leq 1$. Finally, normalize.

c. What are the densities of X and Y when they are $N(0, 1)$?

h.8.2.3. We can make a transform T, U such that $T = X/Y$ and $U = X$ to use a 2D transformation. Compute x and y as functions of t and u . Then the Jacobian.

h.8.2.4. You might want to follow the approach of BH.8.18.

h.8.2.5. Use the bank-post office Story 8.5.1 to see that T and W are independent.

h.8.2.6. a. See BH.8.5.1. The exponential is a special case of the gamma distribution. See also BH.8.34.c. T_1/T_2 is a function of $T_1/(T_1 + T_2)$.

b. This can be solved with a joint distribution function and integration over the event $\{T_1 < T_2\}$. However, we can use Exercise BH.7.10 or BH.7.1.24.

c. First she has to wait for the first server to become free. This is the minimum of the two exponentials. With $P\{T_1 < T_2\}$ server 1 is the first. What is the probability that the other server is empty first? Then, once she is at a server, what is her expected service time? The total time in the system is the time in queue plus the service time.

h.8.2.7. Apply beta-binomial conjugacy.

h.8.2.8. a.

$$P\{X_j \leq c\} = P\{\log U_j \geq -c\} = P\{U_j \geq e^{-c}\} = P\{1 - U_j \leq 1 - e^{-c}\}.$$

What is the distribution of $1 - U_j$?

b. $\log \Pi_{j=1}^n U_j = \sum_{j=1}^n \log U_j = \sum_{j=1}^n (-X_j)$. But $-X_j \sim \text{Exp}(1)$, hence the sum is just a sum of iid Exp rvs. What is the distribution of this sum?

h.8.2.9. Use BH.4.3.9. Then, start with a geometric rv, then extend to a negative binomial rv.

h.9.1.4.**h.9.1.10.**

h.9.1.21. For a smart argument, use the chicken-egg story. Recall that the number of hatched eggs and the number of unhatched eggs are independent (since $N \sim \text{Pois}(\lambda)$); i.e. $N - X$ and X are independent.

h.9.1.22.

h.9.2.1. Then, $E[T|R] = \sum_j E[R_j] I_{R=j}$, where R_j is the time of route j . We know that $E[R_j] = \mu_j$. Now apply the LOTP.

For b. realize that $V[T] = E[T^2] - (E[T])^2$.

h.9.2.2. Use Adam's law to express $E[X_{n+1}]$ in terms of $E[X_n]$, then use recursion.

h.9.2.4. a. Let Y be the amount purchased by the first customer that comes along, let P be the rv that is 1 if the customer does indeed purchase, and 0 otherwise, and let X be the size of the purchase. Why is $Y = XP$? What is $E[P]$? What is $E[Y|P]$? What is $E[Y^2|P]$. You might want to use BH.9.1.

b. Let $N \sim \text{Pois}(8\lambda)$ be the number of customers that pass by. Given $N = n$, what is $E[S|N]$, where $S = \sum_{i=1}^N X_i P_i$ is the total sales. Now use the law of total expectation. What is $V[S|N]$? Use Eve's law to compute $V[S]$. Bigger hint, read Example 9.6.1.

h.9.2.7. a. $N|\lambda \sim \text{Pois}(\lambda)$.

b. Analogous to BH.9.6.1

c. and d. See BH.8.4.5.

h.9.2.10. Refresh your knowledge of the Beta distributions.

a. Since we include the win, the number of games $T|p$ (since we assume p given) must be $\sim \text{FS}(p)$. Hence, $E[T|p] = 1/p$

To get $E[T]$ use Adam's law. Realize that you have to take the integral with respect to p !

b. $1 + E[G]$ is smaller than the expected time as computed in a. Why is this so?

c. The number of wins, conditional on p , out of n is $X|p \sim \text{Bin}(n, p)$. Then use Beta-Binomial conjugacy.

BTW, I find it easier to think about $f(p, X = k)$ instead of $f(p|X = k)$, since on the event $(p, X = k)$.

$$f(p, X = k) \propto p^{a-1} q^{b-1} \binom{n}{k} p^k q^{n-k} \propto p^{a-1+k} q^{b-1+n-k}.$$

Then, as $f(p|X = k) = f(p, X = k)/P\{X = k\} \sim f(p, X = k)$ (because $P\{X = k\}$ is just a constant) we get the same result up to a scaling factor. But we can use the reasoning of BH.8.3.3 to get the correct constant.

h.9.2.11. a. The prior of p is uniform on $[0, 1]$. But this is equal to Beta(1, 1). Now use Beta-Binomial conjugacy.

b. Write $S_n = \sum_{i=1}^n X_i$. What are $P\{X_{n+1} = 1|p\}$ and $P\{S_n = k|p\}$?

h.9.2.12. a. Recall that the uniform distribution on $[0, 1]$ is $\text{Beta}(a, b)$ with $a = b = 1$. I prefer to write $S_n = \sum_{j=1}^n X_j$. First compute $E[S_n | p]$. Then compute $E[E[S_n | p]]$. Note that the outer expectation is an integral with respect to p and the density of $\text{Beta}(1, 1)$.

For the variance, use Eve's law.

b. Use Beta-Binomial conjugacy. Or use the insights of BH.9.56 and BH.9.57.

c. Bayes Billiards.

h.10.1.3.

h.10.1.4.

h.10.1.5.

h.10.1.6.

h.10.1.10. Use that $X \geq X I_{X \geq a} \geq a I_{X \geq a}$.

h.10.1.13. Consider $P\{X \geq 2\}$ and compare Chernoff's inequality to Markov's inequality.

h.10.1.14.

h.10.1.15.

h.10.1.16.

h.10.1.18.

h.10.1.19.

h.10.1.20.

h.10.1.21.

h.10.1.22.

h.10.2.2. First check the assumption that $Y \neq aX$, for some $a > 0$; why is it there? Then, take a suitable g in Jensen's inequality. Bigger hint: $g(x) = 1/x$.

In the solution guide, the authors do not explain the $>$, while in Jensen's inequality there is a \leq . To see why the $>$ is allowed here, rethink the assumption in the exercise, and reread Theorem 10.1.5.

Finally, at what p is $p(1-p)$ maximal?

h.10.2.3. Apply the idea of BH.10.1.3 to $W = (X - \mu)^2$.

h.10.2.4. a. Jensen's inequality, $g(x) = e^x$

b. Use symmetry: X and Y are iid.

c. Which set of events is larger?

d. Use Jensen's inequality and Cauchy-Schwarz.

e. Eve's law.

f. Use Markov's inequality and the triangle inequality

h.10.2.7. The idea is to prove that the MGF of X_n converges to the MGF of a $N(\mu, \sigma^2)$ rv as $n \rightarrow \infty$. Thus, read and follow the proof of the CTL, BH.10.3.1.

What are $E[X_n]$ and $V[X_n]$ if $X \sim \text{Pois}(n)$? Once you know that, explain that the MGF of the standardized version of X_n is equal to $\exp\{-n + s\sqrt{n} + ne^{-s/\sqrt{n}}\}$.

Perhaps you should do BH.10.27 first.

h.10.2.8. a. See BH.10.3.7. Try to convert the recursion for Y_n to a form as in that example.

b. Just substitute α in the relevant formula of part a.

h.10.3.2. Use a substitution. What is the derivative of the survivor function?

h.10.3.3. If you need to prove something for all natural numbers n , it is always good to try using mathematical induction, especially since we already know something relating X_{n+1}^* and X_n^* . Note that you can again use the same substitution as in the previous exercise. After that, you need another substitution.

h.10.3.4. Do you recognize the PDF of X_n^* for $X \sim \text{Exp}(\lambda)$?

h.11.0.1. By the previous exercise, realize that $L \sim \text{Exp}(2\lambda)$. Use these properties.

h.11.0.10. What is the domain of V on each of the intervals $(-3, 0)$ and $[0, 2)$? For the final part, combining the results into one PDF: Use LOTP, conditioning on $U \geq 0$.

h.11.0.19.

h.11.0.22. Use [11.0.21]

h.11.0.23. Use [11.0.22] and the definitions.

h.11.0.26. Use Chebyshev's inequality; then take the limit on both sides.

SOLUTIONS

s.7.1.1. Check the definitions of the book.

Mistake: To say that $P\{X = x\}$ is the PMF for a continuous random variable is wrong, because $P\{X = x\} = 0$ when X is continuous.

Why is $P\{1 < x \leq 4\}$ wrong notation? hint: X should be a capital. What is the difference between X and x ?

s.7.1.2. This example shows why joint distributions are important! Any experiment that involves a sequence of measurements, such as multiple throws of a coin, or the weighing of a bunch of chimpanzees, we have to deal with joint CDFs and PMFs.

s.7.1.3. Here, we deal with two rvs, and we have to specify how they depend. In the present case $P\{X_1 = H, X_2 = H\} = P\{X_1 = H\}$ and $P\{X_1 = T, X_2 = T\} = P\{X_1 = T\}$, $P\{X_1 = H, X_2 = T\} = P\{X_1 = T, X_2 = H\} = 0$. Note that with this, we specified the joint PMF on all possible outcomes.

s.7.1.4.

$$f_X(x) = \int_0^1 f_{X,Y}(x, y) dy = 2 \int_0^1 I_{x \leq y} dy = 2 \int_x^1 dy = 2(1 - x) \quad (13.0.1)$$

$$f_Y(y) = \int_0^1 f_{X,Y}(x, y) dx = 2 \int_0^1 I_{x \leq y} dx = 2 \int_0^y dx = 2y. \quad (13.0.2)$$

But $f_{X,Y}(x, y) \neq f_X(x)f_Y(y)$, hence X, Y are dependent.

$$F_{X,Y}(x, y) = \int_0^x \int_0^y f_{X,Y}(u, v) dv du \quad (13.0.3)$$

$$= 2 \int_0^x \int_0^y I_{u \leq v} dv du \quad (13.0.4)$$

$$= 2 \int_0^x \int_0^y I_{u \leq v} I_{0 \leq v \leq y} dv du \quad (13.0.5)$$

$$= 2 \int_0^x \int_0^y I_{u \leq v \leq y} dv du \quad (13.0.6)$$

$$= 2 \int_0^x [y - u]^+ du, \quad (13.0.7)$$

because $u \geq y \implies I_{u \leq v \leq y} = 0$. Now, if $y > x$,

$$2 \int_0^x [y - u]^+ du = 2 \int_0^x (y - u) du = 2yx - x^2, \quad (13.0.8)$$

while if $y \leq x$,

$$2 \int_0^x [y - u]^+ du = 2 \int_0^y (y - u) du = 2y^2 - y^2 = y^2 \quad (13.0.9)$$

Make a drawing of the support of $f_{X,Y}$ to help to understand this better.

s.7.1.5.

$$\partial_x \partial_y F_{X,Y}(x, y) = \partial_x \partial_y F_X(x) F_Y(y) = \partial_x F_X(x) \partial_y F_Y(y) = f_X(x) f_Y(y).$$

s.7.1.6.

$$\frac{F_{X,Y}(x, y)}{F_X(x)} = \frac{P\{X \leq x, Y \leq y\}}{P\{X \leq x\}} = P\{Y \leq y, X \leq x | X \leq x\} = P\{Y \leq y | X \leq x\}. \quad (13.0.10)$$

It is a big mistake to write $F_{X,Y}(x, y) = P\{X = x, Y = y\}$. If you wrote this, recheck the definitions of BH.

s.7.1.7. $P\{X = 0, Y = 0\} = 1/3 \cdot 3/4$, $P\{X = 0, Y = 1\} = 1/3 \cdot 1/4$, and so on.

If we have one column with $Y = 0$ and the other with $Y = 1$, then the sum over the columns are $P\{Y = 0\}$ and $P\{Y = 1\}$. The row sum for row i are $P\{X = i\}$.

Changing the values will (most of the time) make X and Y dependent. But, what if we changes the values such that $P\{X = 0, Y = 0\} = 1$? Are X and Y then again independent? Check the conditions again.

s.7.1.8. The number of produced items (laid eggs) is N . The probability of hatching is p , that is, an item is ok. The hatched eggs are the good items.

s.7.1.9. For X, Y to be independent, it is necessary that $F_{X,Y}(x, y) = F_X(x)F_Y(y)$ for all x, y , not just one particular choice. (This is an example that satisfying a necessary condition is not necessarily sufficient.)

s.7.1.10. Many answers are possible here, depending on extra assumptions you make. Here is one. Suppose, just by change, the fraction of taller guys in the street is a bit higher than the population fraction. Assuming that taller (shorter) people prefer taller (shorter) spouses, there must be a dependence between the height of the men and the woman. This is because when selecting a man, I can also select his wife.

From this exercise you should memorize that *independence is a property of the joint CDF, not of the rvs.*

Mistake: $P\{Y\}$ is wrong notation wrong because we can only compute the probability of an event, such as $\{Y \leq y\}$. But Y itself is not an event.

s.7.1.11. Only when X, Y are independent.

Mistake: independence of X and Y is not the same as the linear independence. Don't confuse these two types of dependene.

s.7.1.12. Given $N = n$, the random variable X has a certain distribution, here binomial.

s.7.1.13. This claim is incorrect, because X, Y are discrete, hence they have a PMF, not a PDF.

Mistake: Someone said that $\partial_x \partial_y$ is not correct notation; however, it is correct! It's a (much used) abbreviation of the much heavier $\partial^2 / \partial x \partial y$. Next, the derivative of the PMF is not well-defined (at least, not within this course. If you object, ok, but then show that you passed a decent course on measure theory.)

s.7.1.14.

$$\begin{aligned}
 P\{T_1 < T_2\} &= E[I_{T_1 < T_2}] = \int_0^\infty \int_0^\infty I_{t_1 < t_2} f_{T_1, T_2}(t_1, t_2) dt_1 dt_2 \\
 &= \int_0^\infty \int_{t_1}^\infty \lambda_1 e^{-\lambda_1 t_1} \lambda_2 e^{-\lambda_2 t_2} dt_2 dt_1 \\
 &= \int_0^\infty \lambda_1 e^{-\lambda_1 t_1} \lambda_2 \int_{t_1}^\infty e^{-\lambda_2 t_2} dt_2 dt_1 \\
 &= \int_0^\infty \lambda_1 e^{-\lambda_1 t_1} e^{-\lambda_2 t_1} dt_1 \\
 &= \int_0^\infty \lambda_1 e^{-\lambda_1 t_1 - \lambda_2 t_1} dt_1 \\
 &= \frac{\lambda_1}{\lambda_1 + \lambda_2}.
 \end{aligned}$$

s.7.1.15. We have

$$E[(X - Y)^2] = \int_{-\infty}^\infty \int_{-\infty}^\infty (x - y)^2 f_{X, Y}(x, y) dx dy \quad (13.0.11)$$

$$= \int_0^1 \int_0^1 (x - y)^2 dx dy \quad (13.0.12)$$

$$= \int_0^1 \int_0^1 (x^2 - 2xy + y^2) dx dy \quad (13.0.13)$$

$$= \int_0^1 \int_0^1 x^2 dx dy - 2 \int_0^1 \int_0^1 xy dx dy + \int_0^1 \int_0^1 y^2 dx dy \quad (13.0.14)$$

$$= \int_0^1 x^2 dx - 2 \int_0^1 \int_0^1 xy dx dy + \int_0^1 y^2 dy \quad (13.0.15)$$

$$= 1/3 - 2 \cdot 1/2 \cdot 1/2 + 1/3. \quad (13.0.16)$$

s.7.1.16.

$$a < b \implies P\{a < X < b\} = F(b) - F(a) = [F(b) - F(a)]^+ \quad (13.0.17)$$

$$a \geq b \implies P\{a < X < b\} = 0 = [F(b) - F(a)]^+, \quad (13.0.18)$$

where the last equality follows from the fact that F is increasing.

s.7.1.18.

$$\int_{-\infty}^{\infty} I_{0 \leq x \leq 3} dx = \int_0^3 dx = 3.$$

s.7.1.19.

$$\int x I_{0 \leq x \leq 4} dx = \int_0^4 x dx = 16/2 = 8.$$

s.7.1.20.

$$\begin{aligned} \iint xy I_{0 \leq x \leq 3} I_{0 \leq y \leq 4} dx dy &= \int_0^3 x \int_0^4 y dy dx \\ &= \int_0^3 x \frac{y^2}{2} \Big|_0^4 dx \\ &= \int_0^3 x \cdot 8 dx = 8 \cdot 9/2 = 4 \cdot 9. \end{aligned}$$

s.7.1.21. Two solutions. First we integrate over y .

$$\iint I_{0 \leq x \leq 3} I_{0 \leq y \leq 4} I_{x \leq y} dx dy = \int I_{0 \leq x \leq 3} \int I_{0 \leq y \leq 4} I_{x \leq y} dy dx \quad (13.0.19)$$

$$= \int I_{0 \leq x \leq 3} \int I_{\max\{x, 0\} \leq y \leq 4} dy dx \quad (13.0.20)$$

$$= \int_0^3 \int_{\max\{x, 0\}}^4 dy dx \quad (13.0.21)$$

$$= \int_0^3 y \Big|_{\max\{x, 0\}}^4 dx \quad (13.0.22)$$

$$= \int_0^3 (4 - \max\{x, 0\}) dx \quad (13.0.23)$$

$$= 12 - \int_0^3 \max\{x, 0\} dx \quad (13.0.24)$$

$$= 12 - \int_0^3 x dx \quad (13.0.25)$$

$$= 12 - 9/2. \quad (13.0.26)$$

Let's now instead first integrate over x .

$$\iint I_{0 \leq x \leq 3} I_{0 \leq y \leq 4} I_{x \leq y} dx dy = \int I_{0 \leq y \leq 4} \int I_{0 \leq x \leq 3} I_{x \leq y} dx dy \quad (13.0.27)$$

$$= \int_0^4 \int I_{0 \leq x \leq \min\{3, y\}} dx dy \quad (13.0.28)$$

$$= \int_0^4 \int_0^{\min\{3, y\}} dx dy \quad (13.0.29)$$

$$= \int_0^4 \min\{3, y\} dy \quad (13.0.30)$$

$$= \int_0^3 \min\{3, y\} dy + \int_3^4 \min\{3, y\} dy \quad (13.0.31)$$

$$= \int_0^3 y dy + \int_3^4 3 dy \quad (13.0.32)$$

$$= 9/2 + 3. \quad (13.0.33)$$

s.7.1.22. Take c the normalization constant (why is $c = 1/4$), then using the previous exercise

$$P\{Y \leq 2X\} = E[I_{Y \leq 2X}] \quad (13.0.34)$$

$$= c \int_1^3 \int_2^4 I_{y \leq 2x} dy dx \quad (13.0.35)$$

$$= c \int_1^3 \int I_{2 \leq y \leq \min\{4, 2x\}} dy dx \quad (13.0.36)$$

$$= c \int_1^3 [\min\{4, 2x\} - 2]^+ dx \quad (13.0.37)$$

Now make a drawing of the function $[\min\{4, 2x\} - 2]^+$ on the interval $[1, 3]$ to see that

$$\int_1^3 [\min\{4, 2x\} - 2]^+ dx = \int_1^2 (2x - 2) dx + \int_2^3 (4 - 2) dx. \quad (13.0.38)$$

I leave the rest of the computation to you.

s.7.1.23. The covariance might be a large number, which may suggest that the rvs X and Y are ‘very’ dependent. However, when $V[X]$ and $V[Y]$ are also large, the correlation can be small. Thus, correlation is a scaled type of covariance.

s.7.1.24. Answers: no and yes.

We have

$$C = \frac{V[X]}{(E[X])^2}, \quad (13.0.39)$$

which does not equal

$$\text{Corr}(X, X) = \frac{\text{Cov}[X, X]}{\sqrt{V[X]V[X]}} = 1 \quad (13.0.40)$$

in general (for instance, consider a degenerate random variable $X \equiv 1$). Next, consider a $N(1, 100)$ random variable. Then,

$$C = 100/(1^2) = 100 > 1. \quad (13.0.41)$$

s.7.1.25. 1. We have

$$\text{Cov}[X, X] = E[XX] - E[X]E[X] = E[X^2] - E[X]^2 = V[X]. \quad (13.0.42)$$

2. We have

$$\text{Cov}[X, Y] = E[XY] - E[X]E[Y] = E[YX] - E[Y]E[X] = \text{Cov}[Y, X]. \quad (13.0.43)$$

3. We have

$$\text{Cov}[X, c] = E[Xc] - E[X]E[c] = cE[X] - cE[X] = 0. \quad (13.0.44)$$

4. We have

$$\text{Cov}[aX, Y] = E[aXY] - E[aX]E[Y] = a(E[XY] - E[X]E[Y]) = a\text{Cov}[X, Y]. \quad (13.0.45)$$

5. We have

$$\text{Cov}[X + Y, Z] = E[(X + Y)Z] - E[X + Y]E[Z] \quad (13.0.46)$$

$$= E[XZ + YZ] - (E[X] + E[Y])E[Z] \quad (13.0.47)$$

$$= E[XZ] - E[X]E[Z] + E[YZ] - E[Y]E[Z] \quad (13.0.48)$$

$$= \text{Cov}[X, Z] + \text{Cov}[Y, Z]. \quad (13.0.49)$$

s.7.1.26. We have

$$\text{Cov}[X, Y] = E[(X - E[X])(Y - E[Y])] \quad (13.0.50)$$

$$= E[XY - XE[Y] - YE[X] + E[X]E[Y]] \quad (13.0.51)$$

$$= E[XY] - E[X]E[Y] - E[Y]E[X] + E[X]E[Y] \quad (13.0.52)$$

$$= E[XY] - E[X]E[Y]. \quad (13.0.53)$$

When X and Y are independent, then $E[XY] = E[X]E[Y]$, and then $\text{Cov}[X, Y] = 0$.

s.7.1.27. By linearity of the covariance we have

$$\text{Cov}[a(U + V), b(U - V)] = a \left(\text{Cov}[U, b(U - V)] + \text{Cov}[V, b(U - V)] \right) \quad (13.0.54)$$

$$= a \left(b(\text{Cov}[U, U] - \text{Cov}[U, V]) + b(\text{Cov}[V, U] - \text{Cov}[V, V]) \right) \quad (13.0.55)$$

$$= a \left(b(\text{Cov}[U, U] - \text{Cov}[U, V]) + b(\text{Cov}[V, U] - \text{Cov}[V, V]) \right) \quad (13.0.56)$$

$$= ab \left(V[U] - \text{Cov}[U, V] + \text{Cov}[V, U] - V[V] \right) \quad (13.0.57)$$

$$= ab \left(V[U] - V[V] \right). \quad (13.0.58)$$

s.7.1.28. With the hint: $E[XY] = 1/\lambda^2$, when $X, Y \sim \text{Exp}(\lambda)$. Then, $L \sim \text{Exp}(2\lambda)$, since $f_L(x) = 2f_X(x)(1 - F_Y(x)) = 2\lambda e^{-2\lambda x}$. Therefore, $E[L] = 1/2\lambda$. Also, by memoryless, $E[M] = E[L] + E[X] = 3/2\lambda$. Hence, $E[M]E[L] = 3/4\lambda^2$. Hence, $E[ML] - E[M]E[L] = 1/\lambda^2 - 3/4\lambda^2 = 1/4\lambda^2$.

s.7.1.29. We throw 10 fair dice. X_i denotes the number of dice that show the number i , $i = 1, \dots, 6$.

s.7.1.30. No, this does not always hold, see BH.7.5.2. However, it does hold when X and Y are independent.

s.7.1.31. Since X, Y, Z are independent normally distributed variables, (X, Y, Z) is multivariate normally distributed. Hence, every linear combination of X, Y, Z is normally distributed. Note that every linear combination of the elements of W can be written as a linear combination of X, Y, Z . Hence, every linear combination of the elements of W is normally distributed. Hence, W is multivariate normally distributed.

s.7.2.1. Of course $X \in \{0, 1, 2, \dots\}$.

s.7.2.2. For $X > 0$, the first outcome should be a failure. Then, for j failures, we need to fail $j - 1$ times and then once more. For $E[X]$, if there is a success, we don't need another experiment. However, in case of a fail, we need another experiment, and we start again. Thus, $E[X] = q(1 + E[X]) \implies (1 - q)E[X] = q$.

s.7.2.3. With the regular method:

$$\begin{aligned}
 P\{X > j\} &= \sum_{i=j+1}^{\infty} P\{X = i\} \\
 &= p \sum_{i=j+1}^{\infty} q^i \\
 &= p \sum_{i=0}^{\infty} q^{j+1+i} \\
 &= pq^{j+1} \sum_{i=0}^{\infty} q^i \\
 &= pq^{j+1} \frac{1}{1-q} = pq^{j+1} \frac{1}{p} = q^{j+1}.
 \end{aligned}$$

s.7.2.4.

$$\begin{aligned}
E[X] &= \sum_{i=0}^{\infty} i P\{X = i\} \\
&= \sum_{i=0}^{\infty} \sum_{j=0}^{\infty} I_{j < i} P\{X = i\} \\
&= \sum_{j=0}^{\infty} \sum_{i=0}^{\infty} I_{j < i} P\{X = i\} \\
&= \sum_{j=0}^{\infty} \sum_{i=j+1}^{\infty} P\{X = i\} \\
&= \sum_{j=0}^{\infty} P\{X > j\} \\
&= \sum_{j=0}^{\infty} q^{j+1} \\
&= q \sum_{j=0}^{\infty} q^j \\
&= q/(1 - q) = q/p.
\end{aligned}$$

s.7.2.5.

$$\begin{aligned}
P\{X \geq n + m \mid X \geq m\} &= \frac{P\{X \geq n + m, X \geq m\}}{P\{X \geq m\}} \\
&= \frac{P\{X \geq n + m\}}{P\{X \geq m\}} \\
&= \frac{q^{n+m}}{q^m} \\
&= q^n = P\{X \geq n\}.
\end{aligned}$$

s.7.2.6.

$$\begin{aligned}
P\{L \geq k\} &= \sum_i \sum_j I_{\min\{i, j\} \geq k} P\{X = i, Y = j\} \\
&= \sum_{i \geq k} \sum_{j \geq k} P\{X = i\} P\{Y = j\} \\
&= P\{X \geq k\} P\{Y \geq k\} = q^k q^k = q^{2k}.
\end{aligned}$$

$P\{L > i\}$ has the same form as $P\{X > i\}$, but now with q^{2i} rather than q^i .

s.7.2.7. Immediate from the hint and [7.2.4].

s.7.2.8.

$$\begin{aligned}
E[L] + E[X] &= \frac{q^2}{1-q^2} + \frac{q}{1-q} \\
&= \frac{q}{1-q} \left(\frac{q}{1+q} + 1 \right) \\
&= \frac{q}{1-q} \frac{1+2q}{1+q}
\end{aligned}$$

s.7.2.9.

$$\begin{aligned}
2E[X] - E[L] &= 2 \frac{q}{1-q} - \frac{q^2}{1-q^2} \\
&= \frac{q}{1-q} \left(2 - \frac{q}{1+q} \right) \\
&= \frac{q}{1-q} \left(\frac{2+2q}{1+q} - \frac{q}{1+q} \right) \\
&= \frac{q}{1-q} \frac{2+q}{1+q}.
\end{aligned}$$

s.7.2.10.

$$\begin{aligned}
P\{M = k\} &= P\{\max\{X, Y\} = k\} \\
&= p^2 \sum_{i,j} I_{\max\{i,j\}=k} q^i q^j \\
&= 2p^2 \sum_{i,j} I_{i=k} I_{j < k} q^i q^j + p^2 \sum_{i,j} I_{i=j=k} q^i q^j \\
&= 2p^2 q^k \sum_{j < k} q^j + p^2 q^{2k} \\
&= 2p^2 q^k \frac{1-q^k}{1-q} + p^2 q^{2k}
\end{aligned}$$

s.7.2.11. It's just algebra

$$\begin{aligned}
P\{M = k\} &= 2p^2 q^k \frac{1-q^k}{1-q} + p^2 q^{2k} \\
&= 2p q^k (1-q^k) + p^2 q^{2k} \\
&= 2p q^k + (p^2 - 2p) q^{2k} \\
&= 2P\{X = k\} - P\{L = k\},
\end{aligned}$$

where I use that $p^2 - 2p = p(p - 2) = (1 - q)(1 - q - 2) = -(1 - q)(1 + q) = -(1 - q^2)$.

s.7.2.12.

$$\begin{aligned}
E[M] &= \sum_k k P\{M = k\} \\
&= \sum_k k(2P\{X = k\} - P\{L = k\}) \\
&= 2E[X] - E[L].
\end{aligned}$$

s.7.2.13.

$$\begin{aligned}
P\{M \leq k\} &= P\{X \leq k, Y \leq k\} \\
&= P\{X \leq k\} P\{Y \leq k\} \\
&= (1 - P\{X > k\})(1 - P\{Y > k\}) \\
&= (1 - q^{k+1})^2.
\end{aligned}$$

s.7.2.14.

$$\begin{aligned}
P\{M = k\} &= P\{M \leq k\} - P\{M \leq k-1\} \\
&= 1 - 2q^{k+1} + q^{2k+2} - (1 - 2q^k + q^{2k}) \\
&= 2q^k(1 - q) + q^{2k}(q^2 - 1) \\
&= 2P\{X = k\} - q^{2k}(1 - q^2) \\
&= 2P\{X = k\} - P\{L = k\}.
\end{aligned}$$

s.7.2.15.

$$\begin{aligned}
P\{L = i, M = k\} &= 2P\{X = i, Y = k\} I_{k>i} + P\{X = Y = i\} I_{i=k} \\
&= 2p^2 q^{i+k} I_{k>i} + p^2 q^{2i} I_{i=k}.
\end{aligned}$$

s.7.2.16.

$$\begin{aligned}
P\{M = k\} &= \sum_i P\{L = i, M = k\} \\
&= \sum_i (2p^2 q^{i+k} I_{k>i} + p^2 q^{2i} I_{i=k}) \\
&= 2p^2 q^k \sum_{i=0}^{k-1} q^i + p^2 q^{2k} \\
&= 2pq^k(1 - q^k) + p^2 q^{2k} \\
&= 2pq^k + (p^2 - 2p)q^{2k},
\end{aligned}$$

s.7.2.17.

$$\begin{aligned}
 P\{L = i\} &= \sum_k P\{L = i, M = k\} \\
 &= \sum_k (2p^2 q^{i+k} I_{k>i} + p^2 q^{2i} I_{i=k}) \\
 &= 2p^2 q^i \sum_{k=i+1}^{\infty} q^k + p^2 q^{2i} \\
 &= 2p^2 q^{2i+1} \sum_{k=0}^{\infty} q^k + p^2 q^{2i} \\
 &= 2p q^{2i+1} + p^2 q^{2i} \\
 &= p q^{2i} (2q + p) \\
 &= (1 - q) q^{2i} (q + 1), \quad p + q = 1, \\
 &= (1 - q^2) q^{2i}.
 \end{aligned}$$

s.7.2.18. With $t \geq s$,

$$P\{X > t | X > s\} = \frac{P\{X > t\}}{P\{X > s\}} = e^{-\lambda t} e^{\lambda s} = e^{-\lambda(t-s)} = P\{X > t - s\}.$$

s.7.2.19. It is essential that you know both methods to solve this integral.

$$\begin{aligned}
 E[X] &= \int_0^{\infty} x f_X(x) dx \\
 &= \int_0^{\infty} x \lambda e^{-\lambda x} dx \\
 &= -x e^{-\lambda x} \Big|_0^{\infty} + \int_0^{\infty} e^{-\lambda x} dx \\
 &= -\frac{1}{\lambda} e^{-\lambda x} \Big|_0^{\infty} = \frac{1}{\lambda}.
 \end{aligned}$$

Substitution is also a very important technique to solve such integrals. Here we go again:

$$\begin{aligned}
 E[X] &= \int_0^{\infty} x f_X(x) dx \\
 &= \int_0^{\infty} x \lambda e^{-\lambda x} dx \\
 &= \frac{1}{\lambda} \int_0^{\infty} u e^{-u} du,
 \end{aligned}$$

by the substitution $u = \lambda x \implies du = d(\lambda x) \implies du = \lambda dx \implies dx = du/\lambda$. With partial integration (do it!), the integral evaluates to 1.

s.7.2.20. With the hint,

$$G_L(x) = P\{L > x\} = P\{X > x, Y > x\} = G_X(x)^2 = e^{-2\lambda x}.$$

The result follows since $F_L(x) = 1 - G_L(x)$.

s.7.2.21.

$$\begin{aligned}
P\{M \leq v\} &= E[I_{M \leq v}] \\
&= \int_0^\infty \int_0^\infty I_{x \leq v, y \leq v} f_{XY}(x, y) dx dy \\
&= \int_0^\infty \int_0^\infty I_{x \leq v, y \leq v} f_X(x) f_Y(y) dx dy \\
&= \int_0^v f_X(x) dx \int_0^v f_Y(y) dy \\
&= F_X(v) F_Y(v) = (F_X(v))^2.
\end{aligned}$$

s.7.2.22.

$$f_M(v) = F'_M(v) = 2F_X(v)f_X(v) = 2(1 - e^{-\lambda v})\lambda e^{-\lambda v}.$$

s.7.2.23.

$$\begin{aligned}
E[M] &= \int_0^\infty v f_M(v) dv = \\
&= 2\lambda \int_0^\infty v(1 - e^{-\lambda v})e^{-\lambda v} dv = \\
&= 2\lambda \int_0^\infty v e^{-\lambda v} dv - 2\lambda \int_0^\infty v e^{-2\lambda v} dv \\
&= 2E[X] - E[L],
\end{aligned}$$

where the last equality follows from the previous exercises.

s.7.2.24. First the joint distribution. With $u \leq v$,

$$\begin{aligned}
F_{L,M}(u, v) &= P\{L \leq u, M \leq v\} \\
&= 2 \iint I_{x \leq u, y \leq v, x \leq y} f_{X,Y}(x, y) dx dy \\
&= 2 \int_0^u \int_x^v f_Y(y) dy f_X(x) dx && \text{independence} \\
&= 2 \int_0^u (F_Y(v) - F_Y(x)) f_X(x) dx.
\end{aligned}$$

s.7.2.25. Taking partial derivatives,

$$\begin{aligned}
f_{L,M}(u, v) &= \partial_v \partial_u F_{L,M}(u, v) \\
&= 2 \partial_v \partial_u \int_0^u (F_Y(v) - F_Y(x)) f_X(x) dx \\
&= 2 \partial_v \{(F_Y(v) - F_Y(u)) f_X(u)\} \\
&= 2 f_X(u) \partial_v F_Y(v) \\
&= 2 f_X(u) f_Y(v).
\end{aligned}$$

s.7.2.26.

$$\begin{aligned}
f_M(v) &= \int_0^\infty f_{L,M}(u, v) \, du \\
&= 2 \int_0^v f_X(u) f_Y(v) \, du \\
&= 2f_Y(v) \int_0^v f_X(u) \, du \\
&= 2f_Y(v) F_X(v), \\
f_L(u) &= \int_0^\infty f_{L,M}(u, v) \, dv \\
&= 2f_X(u) \int_u^\infty f_Y(v) \, dv \\
&= 2f_X(u) G_Y(u).
\end{aligned}$$

s.7.2.27. First,

$$\begin{aligned}
P\{X/n \approx x\} &= P\{X/n \in [i/n, (i+1)/n]\} = P\{X \in [i, i+1]\} = pq^i \\
&\approx pq^{nx} = \frac{\lambda}{n} \left(1 - \frac{\lambda}{n}\right)^{xn},
\end{aligned}$$

since $p = \lambda/n$, $q = 1 - p = 1 - \lambda/n$, and $i = nx$.

s.7.2.28.

$$M_{X/n}(s) = E\left[e^{sX/n}\right] = \sum_i e^{si/n} pq^i = p \sum_i (qe^{s/n})^i = \frac{p}{1 - qe^{s/n}}.$$

With $p = \lambda/n$ this becomes

$$\begin{aligned}
M_{X/n}(s) &= \frac{\lambda/n}{1 - (1 - \lambda/n)(1 + s/n + 1/n^2 \times (\dots))} \\
&= \frac{\lambda/n}{\lambda/n - s/n + 1/n^2 \times (\dots)} \\
&= \frac{\lambda}{\lambda - s + 1/n \times (\dots)} \\
&\rightarrow \frac{\lambda}{\lambda - s}, \quad \text{as } n \rightarrow \infty,
\end{aligned}$$

where we write $1/n^2 \times (\dots)$ for all terms that will disappear when we take the limit $n \rightarrow \infty$. This is just handy notation to hide details in which we are not interested.

s.7.2.29.

$$\begin{aligned}
E[L/n] &= \frac{1}{n} E[L] = \frac{1}{n} \frac{q^2}{1-q^2} \\
&= \frac{1}{n} \frac{(1-\lambda/n)^2}{1-(1-\lambda/n)^2} \\
&= \frac{1}{n} \frac{1-2\lambda/n+(\lambda/n)^2}{2\lambda/n+(\lambda/n)^2} \\
&= \frac{1-2\lambda/n+(\lambda/n)^2}{2\lambda+\lambda^2/n} \\
&\rightarrow \frac{1}{2\lambda}.
\end{aligned}$$

s.7.2.30. Take $p = \lambda/n$, $q = 1 - \lambda/n$, and $k \approx xn$, hence $k/n \approx x$. Then,

$$\begin{aligned}
P\{M/n = k/n\} &= 2pq^{k/n}(1-q^{k/n}) + p^2q^{2k/n} \\
&= 2\frac{\lambda}{n}\left(1-\frac{\lambda}{n}\right)^{k/n}\left(1-\left(1-\frac{\lambda}{n}\right)^{k/n}\right) + \frac{\lambda^2}{n^2}\left(1-\frac{\lambda}{n}\right)^{2k/n} \\
&\rightarrow 2\lambda dx e^{-\lambda x}\left(1-e^{-\lambda x}\right) + \lambda^2 dx^2 e^{-2\lambda}.
\end{aligned}$$

Now observe that the second term, proportional to dx^2 can be neglected.

s.7.2.31.

$$\begin{aligned}
F_{L,M-L}(x,y) &= P\{L \leq x, M-L \leq y\} \\
&= 2P\{X \leq x, Y-X \leq y, X \leq Y\} \\
&= 2 \int_0^\infty \int_0^\infty I_{u \leq x, v-u \leq y, u \leq v} f_{X,Y}(u,v) du dv \\
&= 2 \int_0^\infty \int_0^\infty I_{u \leq x, v-u \leq y, u \leq v} \lambda^2 e^{-\lambda u} e^{-\lambda v} du dv \\
&= 2 \int_0^x \int_0^\infty I_{u \leq v \leq u+y} \lambda^2 e^{-\lambda u} e^{-\lambda v} dv du \\
&= 2 \int_0^x \lambda e^{-\lambda u} \int_u^{u+y} \lambda e^{-\lambda v} dv du \\
&= 2 \int_0^x \lambda e^{-\lambda u} (-e^{-\lambda v}) \Big|_u^{u+y} du \\
&= 2 \int_0^x \lambda e^{-\lambda u} (e^{-\lambda u} - e^{-\lambda(u+y)}) du \\
&= 2\lambda \int_0^x e^{-2\lambda u} du - 2\lambda \int_0^x e^{-\lambda(2u+y)} du \\
&= 2\lambda \int_0^x e^{-2\lambda u} du - 2\lambda e^{-\lambda y} \int_0^x e^{-2\lambda u} du \\
&= (1 - e^{-\lambda y}) 2\lambda \int_0^x e^{-2\lambda u} du \\
&= (1 - e^{-\lambda y}) (-e^{-2\lambda u}) \Big|_0^x \\
&= (1 - e^{-\lambda y})(1 - e^{-2\lambda x}).
\end{aligned}$$

s.7.2.32. As $F_{L,M-L}(x,y) = F_Y(y)F_L(x)$. So the CDF factors as a function of x only and a function of y only. This implies that L and $M-L$ are independent, and moreover that $F_{M-L}(y) = F_Y(y)$, so $M-L \sim Y$. We can also see this from the joint PDF:

$$f_{L,M-L}(x,y) = \partial_x \partial_y (F_Y(y)F_L(x)) = f_Y(y)f_L(x),$$

so the joint PDF (of course) also factors. The independence now follows from BH 7.1.21. Because L and $M-L$ are independent, the conditional density equals the marginal density:

$$f_{M-L|L}(y|x) = \frac{f_{L,M-L}(x,y)}{f_L(x)} = \frac{f_Y(y)f_L(x)}{f_L(x)} = f_Y(y).$$

s.7.2.34. Suppose $j > 0$ (for $j = 0$ the maths is the same). Then,

$$P\{M-L = j\} = 2 \sum_{i=0}^{\infty} P\{X = i, Y = i+j\} = 2 \sum_{i=0}^{\infty} p q^i p q^{i+j} = 2 p^2 q^j \sum_{i=0}^{\infty} q^{2i} = 2 p^2 q^j / (1 - q^2).$$

s.7.2.35. Because either $M = L$ or $M > L$. Recall from earlier work that the factor 2 in the second equality follows from the fact that X, Y iid.

s.7.2.36.

$$\begin{aligned}
E[(Y - X)I_{Y>X}] &= p^2 \sum_{ij} (j - i) I_{j>i} q^i q^j \\
&= p^2 \sum_i q^i \sum_{j=i+1}^{\infty} (j - i) q^j \\
&= p^2 \sum_i q^i q^i \sum_{k=1}^{\infty} k q^k, \quad k = j - i \\
&= p \sum_i q^{2i} E[X] \\
&= \frac{p}{1 - q^2} E[X] \\
&= \frac{p}{1 - q^2} \frac{q}{p} \\
&= \frac{q}{1 - q^2}.
\end{aligned}$$

s.7.2.37.

$$E[L] + 2E[(Y - X)I_{Y>X}] = \frac{q^2}{1 - q^2} + \frac{2q}{1 - q^2} = \frac{q}{1 - q} \frac{q + 2}{1 + q},$$

where I use that $1 - q^2 = (1 - q)(1 + q)$.

s.7.2.38. To see why this might be true, I reason like this. After ‘seeing’ L , we want to restart. Let Z be the time from the restart to M . When $Z \sim \text{Geo}(p)$, it might happen that $Z = 0$ (with positive probability p). But if $Z = 0$, then $M = L$, and in that case, we should not restart. Hence, if $Z \sim \text{Geo}(p)$ we are ‘double counting’ when $Z = 0$. By including the condition $M > L$ and by taking $Z \sim \text{FS}(p)$ (so that $Z > 0$) I can prevent this.

s.7.2.39. With the hint:

$$E[Z I_{M>L}] = E[Z] E[I_{M>L}] = E[Z] P\{M > L\} = \frac{1}{p} \frac{2pq}{1 - q^2} = \frac{2q}{1 - q^2},$$

We know that $E[Z] = 1 + E[X] = 1 + q/p = 1/p$, while

$$\begin{aligned}
P\{M > L\} &= 1 - P\{X = Y\} = 1 - \sum_{i=0}^{\infty} P\{X = Y = i\} \\
&= 1 - \frac{p^2}{1 - q^2} = \frac{1 - q^2 + p^2}{1 - q^2} = \frac{2pq}{1 - q^2}.
\end{aligned}$$

s.7.2.40. By independence,

$$E[Z/n I_{M>L}] = E[Z/n] P\{M > L\}.$$

Then,

$$P\{M > L\} = \frac{2pq}{1 - q^2} = \frac{2\lambda/n(1 - \lambda/n)}{1 - (1 - \lambda/n)^2} = \frac{2\lambda/n(1 - \lambda/n)}{2\lambda/n - \lambda^2/n^2} = \frac{2(1 - \lambda/n)}{2 - \lambda/n} \rightarrow 1,$$

and

$$E[Z/n] = \frac{1}{n} E[Z] = \frac{1}{n} \frac{1}{p} = \frac{1}{n\lambda/n} = 1/\lambda.$$

s.7.3.1. a.

$$F(2, 5) = P(X \leq 2, Y \leq 5) = P(X \leq 2, Y \leq 4) = F(2, 4) = \frac{1}{4}$$

The second step is valid since the cumulative distribution function does not change by changing y from 5 to 4 thanks to property 3.

b. To obtain the joint pdf, use that $f_{X,Y}(x, y) = \frac{\partial^2}{\partial y \partial x} F(x, y) = \frac{\partial}{\partial y} \left(\frac{\partial}{\partial x} F(x, y) \right)$.

Since $\frac{\partial}{\partial x} F(x, y) = \frac{1}{4}(x-1)(y-2)$ for $1 < x < 3$, and $\frac{\partial}{\partial x} F(x, y) = 0$ for other values of x , we have that

$$f_{X,Y}(x, y) = \begin{cases} \frac{1}{4}(x-1), & \text{for } 1 < x < 3 \text{ and } 2 < y < 4, \\ 0, & \text{elsewhere.} \end{cases}$$

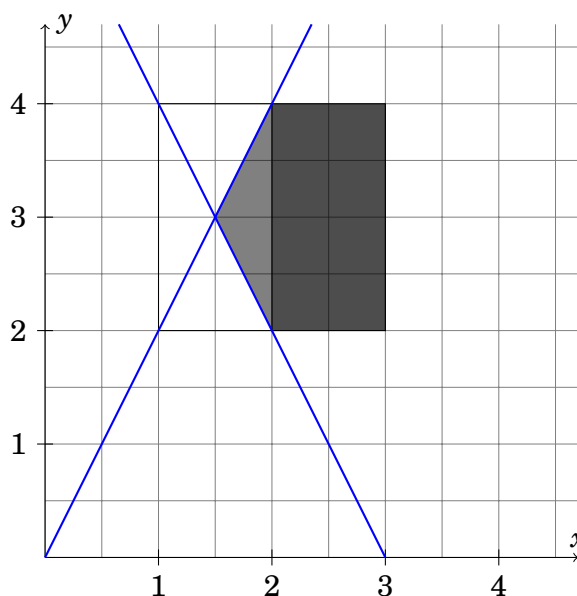
c. The simplest way of solving this question is by writing

$$\begin{aligned} P(2 < X < 3, 2 < Y < 4) &= F(3, 4) - F(2, 4) - F(3, 2) + F(2, 2) \\ &= 1 - \frac{1}{4} - 0 + 0 = \frac{3}{4}. \end{aligned}$$

Alternatively, one can integrate over the pdf from (b) to obtain the same result:

$$\begin{aligned} P(2 < X < 3, 2 < Y < 4) &= \int_2^3 \int_2^4 f(x, y) dy dx \\ &= \frac{1}{4} \int_2^3 (x-1) y \Big|_{y=2}^{y=4} dx \\ &= \frac{1}{4} (x^2 - 2x) \Big|_2^3 = \frac{1}{4} (9 - 6 - 4 + 4) = \frac{3}{4}. \end{aligned}$$

d. First, draw the integration area:



The domain on which the density is non-zero, is the complete shaded area. The downward-sloping line represents $2x + y = 6$ and the upward-sloping line is $y = 2x$.

We already know the integral over the dark shaded area from the previous subquestion. What remains is the lighter shaded triangular part on the left side.

First, we need to calculate the intersection of the two curves, which can be found by solving $6 - 2x = 2x$, which gives $x = \frac{3}{2}$, and consequently $y = 3$.

The integral of the triangular part of the dark shaded region is

$$\begin{aligned} \int_{3/2}^2 \int_{6-2x}^{2x} \frac{1}{4}(x-1)dydx &= \int_{3/2}^2 \frac{1}{4}(x-1)y \Big|_{6-2x}^{2x} dx \\ &= \frac{1}{4} \int_{3/2}^2 (x-1)2x - (x-1)(6-2x)dx \\ &= \frac{1}{4} \int_{3/2}^2 (x-1)(4x-6)dx \\ &= \frac{1}{4} \left[\frac{4}{3}x^3 - 5x^2 + 6x \right]_{3/2}^2 \\ &= \frac{5}{48} \approx 0.1042 \end{aligned}$$

Finally, the joint probability asked for in the question is given by

$$P(Y < 2X, 2X + Y > 6) = \frac{5}{48} + \frac{3}{4} = \frac{41}{48} \approx 0.8542$$

s.7.3.2. For $r \geq 0$, we have

$$\begin{aligned} F_R(r) &= P(R \leq r) \\ &= P(X \leq Vr) \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{vr} f_{X,V}(x,v) dx dv \\ &= \lambda \mu \int_0^{\infty} \int_0^{vr} e^{-\lambda x} e^{-\mu v} dx dv \\ &= \lambda \mu \int_0^{\infty} e^{-\mu v} \left[-\frac{1}{\lambda} e^{-\lambda x} \right]_0^{vr} dv \\ &= -\mu \int_0^{\infty} e^{-\mu v} (e^{-\lambda vr} - 1) dv \\ &= -\mu \left[-\frac{1}{\mu + \lambda r} e^{-(\mu + \lambda r)v} + \frac{1}{\mu} e^{-\mu v} \right]_0^{\infty} \\ &= -\mu \left[\frac{1}{\mu + \lambda r} - \frac{1}{\mu} \right] \\ &= \frac{\lambda r}{\mu + \lambda r}, \end{aligned}$$

while $F_R(r) = 0$ when $r < 0$ since both X and V are nonnegative.

We see that (1) $F_R(-\infty) = 0$, (2) $F_R(\infty) = 1$, and (3) $F_R(r)$ is monotonically increasing in r , so $F_R(r)$ satisfies the conditions for being a valid CDF.

s.7.3.3. a. We integrate $f_{X,Y}(x,y)$ over its domain

$$\begin{aligned}
 & \int_0^{1/2} \int_0^x cxydydx + \int_{1/2}^1 \int_0^{1-x} cxydydx \\
 &= c \left[\frac{1}{2} \int_0^{1/2} x^3 dx + \frac{1}{2} \int_{1/2}^1 x(1-x)^2 dx \right] \\
 &= \frac{1}{2} c \left\{ \left[\frac{1}{4} x^4 \right]_0^{1/2} + \left[\frac{1}{2} x^2 - \frac{2}{3} x^3 + \frac{1}{4} x^4 \right]_{1/2}^1 \right\} \\
 &= \frac{1}{2} c \left\{ \frac{1}{4} \frac{1}{16} + \frac{1}{2} - \frac{2}{3} + \frac{1}{4} - \frac{1}{2} \frac{1}{4} + \frac{2}{3} \frac{1}{8} - \frac{1}{4} \frac{1}{16} \right\} \\
 &= \frac{1}{2} c \frac{12 - 16 + 6 - 3 + 2}{24} \\
 &= \frac{c}{48}
 \end{aligned}$$

Since this integral should equal 1, $c = 48$.

Alternative Rewrite the probability density function to

$$f_{X,Y}(x,y) = \begin{cases} cxy & 0 \leq y \leq \frac{1}{2}, \quad y \leq x \leq 1-y \\ 0 & \text{otherwise} \end{cases}$$

Then,

$$\begin{aligned}
 \int_0^{1/2} \int_y^{1-y} cxydx dy &= c \int_0^{1/2} y \frac{1}{2} ((1-y)^2 - y^2) dy \\
 &= \frac{1}{2} c \int_0^{1/2} (y - 2y^2) dy \\
 &= \frac{1}{2} c \left[\frac{1}{2} y^2 - \frac{2}{3} y^3 \right]_0^{1/2} \\
 &= \frac{1}{2} c \left[\frac{1}{8} - \frac{2}{3} \frac{1}{8} \right] = \frac{c}{48}
 \end{aligned}$$

Since this integral should equal 1, $c = 48$.

b. The conditional density function is given by

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)}$$

We first need the marginal density of Y .

$$\begin{aligned}
 f_Y(y) &= 48 \int_y^{1-y} xy dx \\
 &= 48y \left[\frac{1}{2}(1-y)^2 - \frac{1}{2}y^2 \right] \\
 &= 24y [1 - 2y + y^2 - y^2] \\
 &= 24y(1 - 2y) \quad \text{for } 0 \leq y \leq \frac{1}{2}
 \end{aligned}$$

and $f_Y(y) = 0$ otherwise.

Not required: We can check that this is a valid density function:

$$\begin{aligned}
 \int_0^{1/2} 24y(1 - 2y) dy &= 24 \left[\frac{1}{2}y^2 - \frac{2}{3}y^3 \right]_0^{1/2} \\
 &= 24 \left[\frac{1}{2} \cdot \frac{1}{4} - \frac{2}{3} \cdot \frac{1}{8} \right]_0^{1/2} \\
 &= 1
 \end{aligned}$$

b. Now we can obtain the conditional density function

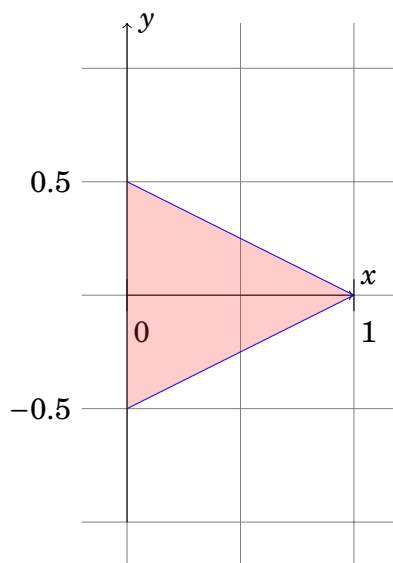
$$\begin{aligned}
 f_{X|Y}(x|y) &= \frac{f_{X,Y}(x,y)}{f_Y(y)} \\
 &= \frac{48xy}{24y(1 - 2y)} = \frac{2x}{1 - 2y} \quad \text{for } y \leq x \leq 1 - y, \text{ and } 0 \leq y < \frac{1}{2}
 \end{aligned}$$

and $f_{X|Y}(x|y) = 0$ otherwise.

This is a valid density function since $f_{X|Y}(x|y) \geq 0$, and

$$\begin{aligned}
 \int_y^{1-y} f_{X|Y}(x|y) dx &= \frac{2}{1 - 2y} \left[\frac{1}{2}x^2 \right]_y^{1-y} \\
 &= 1
 \end{aligned}$$

s.7.3.4. a. The joint PDF is nonzero above the red-shaded area in the following graph. (draw, draw, draw!)



For the PDF for Y , using the graph above,

$$\begin{aligned}
 f_Y(y) &= \int_{-\infty}^{\infty} f_{X,Y}(x,y) dx \\
 &= \int_0^{1-2|y|} 2 dx \\
 &= 2(1-2|y|) \quad \text{for } -1/2 < y < 1/2,
 \end{aligned}$$

and 0 otherwise. Since $-1/2 < y < 1/2$, we have $f_Y(y) \geq 0$. Also,

$$\begin{aligned}
 \int_{-\infty}^{\infty} f_Y(y) dy &= \int_{-1/2}^{1/2} 2(1-2|y|) dy \\
 &= 2 \int_0^{1/2} (1-2y) dy + 2 \int_{-1/2}^0 (1+2y) dy \\
 &= 2 \left(\frac{1}{2} - \frac{1}{4} + \frac{1}{2} - \frac{1}{4} \right) = 1.
 \end{aligned}$$

Probability density function for X .

$$\begin{aligned}
 f_X(x) &= \int_{-\infty}^{\infty} f_{X,Y} dy \\
 &= \int_{-\frac{1}{2}(1-x)}^{\frac{1}{2}(1-x)} 2 dy \\
 &= 2(1-x) \quad \text{for } 0 \leq x < 1,
 \end{aligned}$$

and 0 otherwise. We have $f_X(x) \geq 0$, and also

$$\begin{aligned}
 \int_{-\infty}^{\infty} f_X(x) dx &= \int_0^1 2(1-x) dx \\
 &= 2 \left[x - \frac{1}{2}x^2 \right]_0^1 \\
 &= 2 \left(1 - \frac{1}{2} \right) = 1.
 \end{aligned}$$

Since $-1/2 < y < 1/2$, we have $f_Y(y) \geq 0$. Also,

$$\begin{aligned}\int_{-\infty}^{\infty} f_Y(y) dy &= \int_{-1/2}^{1/2} 2(1 - 2|y|) dy \\ &= 2 \int_0^{1/2} (1 - 2y) dy + 2 \int_{-1/2}^0 (1 + 2y) dy \\ &= 2 \left(\frac{1}{2} - \frac{1}{4} + \frac{1}{2} - \frac{1}{4} \right) = 1.\end{aligned}$$

Probability density function for X .

$$\begin{aligned}f_X(x) &= \int_{-\infty}^{\infty} f_{X,Y} dy \\ &= \int_{-\frac{1}{2}(1+x)}^{\frac{1}{2}(1+x)} 2 dy \\ &= 2(1+x) \quad \text{for } -1 < x \leq 0,\end{aligned}$$

and 0 otherwise. We have $f_X(x) \geq 0$, and also

$$\begin{aligned}\int_{-\infty}^{\infty} f_X(x) dx &= \int_{-1}^0 2(1+x) dx \\ &= 2 \left[x + \frac{1}{2}x^2 \right]_{-1}^0 \\ &= 2 \left(1 - \frac{1}{2} \right) = 1.\end{aligned}$$

b. Using the definition of a conditional probability density function, we have

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)} = \frac{1}{1-2|y|} \quad \text{for } -\frac{1}{2} < y < \frac{1}{2}, 0 < x < 1-2|y|$$

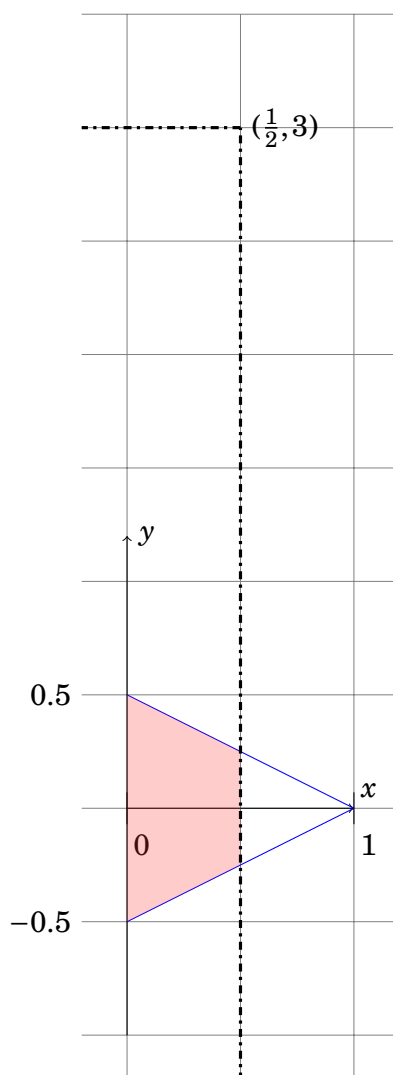
and 0 otherwise.

The required expected value is (bounds are crucial!)

$$E[X|Y=y] = \int_0^{1-2|y|} \frac{x}{1-2|y|} dx = \frac{1}{2}(1-2|y|).$$

If $y = -1/2$ or $y = 1/2$, we find $E[X|Y=y] = 0$, which makes sense based on the figure above. Similarly, if $y = 0$, we find $E[X|Y=y] = 1/2$ since the density of x is uniform over $[0, 1]$.

The point $(\frac{1}{2}, 3)$ is located as in the picture below. We need to integrate $f_{X,Y}(x,y)$ over the entire area on the lower left side of this point. It is crucial to realize that the joint PDF is only nonzero in the red shaded area. Moreover, it's very helpful that the PDF has a constant value of 2 above this area. So we just need to calculate the area of the red shape in the following graph and multiply this by 2.



The dash dotted line intersects the line $y = \frac{1}{2}(1 - x)$ at $x = \frac{1}{2}$, so $y = \frac{1}{4}$. The line intersects $y = -\frac{1}{2}(1 - x)$ at $x = \frac{1}{2}$, so $y = -\frac{1}{4}$. The whole area of the triangle is $\frac{1}{2}$. The area *not* in red is $\frac{1}{2} \cdot (\frac{1}{4} - (-\frac{1}{4})) \cdot (1 - \frac{1}{2}) = \frac{1}{8}$. So

$$F_{X,Y}\left(\frac{1}{2}, 3\right) = 2\left(\frac{1}{2} - \frac{1}{8}\right) = \frac{3}{4}.$$

s.7.3.5. a. $f_{X,Y}(x, y)$ is a joint probability density function if

1. If $f_{X,Y}(x, y)$ satisfies

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x, y) dx dy = 1$$

2. $f_{X,Y}(x, y) \geq 0$ for all x and y .

We have

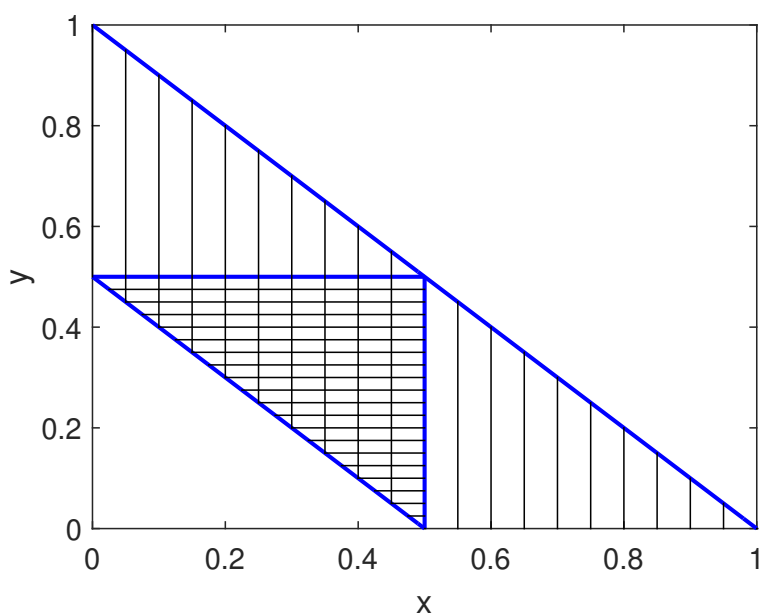
$$\begin{aligned}\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{X,Y}(x,y) dy dx &= \int_0^1 \int_0^{1-x} \frac{c}{1-x} dy dx \\ &= c \int_0^1 \frac{1-x}{1-x} dx = 1\end{aligned}$$

So to satisfy condition (1), we need to set $c = 1$.

Check Condition (2): $f_{X,Y}(x,y) \geq 0$ for all x,y .

For $0 < x < 1$, $f_{X,Y}(x,y) > 0$. Outside of this interval $f_{X,Y}(x,y) = 0$. So, we have that $f_{X,Y}(x,y) \geq 0$ for all x,y .

b. The following graph is used for questions b and c .



To obtain $P(X + Y > 1/2)$, we need to integrate $f_{X,Y}$ over the vertically hatched area. We do this by integrating over the larger triangle defined by $x + y < 1$, and then subtract the white triangle in the lower left corner defined by $x + y < \frac{1}{2}$.

$$\begin{aligned}
 P\left(X + Y > \frac{1}{2}\right) &= \int_0^1 \int_0^{1-x} \frac{1}{1-x} dy dx - \int_0^{1/2} \int_0^{1/2-x} \frac{1}{1-x} dy dx \\
 &= 1 - \int_0^{1/2} \frac{1/2-x}{1-x} dx \\
 &= 1 - \int_0^{1/2} \left(1 - \frac{1}{2} \frac{1}{1-x}\right) dx \\
 &= 1 - \frac{1}{2} + \frac{1}{2} \int_0^{1/2} \frac{1}{1-x} dx \\
 &= \frac{1}{2} + \frac{1}{2} [-\ln(1-x)]_0^{1/2} \\
 &= \frac{1}{2} - \frac{1}{2} \ln\left(\frac{1}{2}\right) = 0.8466
 \end{aligned}$$

c. Using the definition of conditional probability

$$P\left(X < \frac{1}{2}, Y < \frac{1}{2} \mid X + Y > \frac{1}{2}\right) = \frac{P\left(X < \frac{1}{2}, Y < \frac{1}{2}, X + Y > \frac{1}{2}\right)}{P\left(X + Y > \frac{1}{2}\right)}$$

The integral in the numerator is the integral over the horizontally hatched triangle. Easiest is to first integrate over the square $[0, 1/2] \times [0, 1/2]$ and then subtract the lower white triangle, which we have already calculated in the previous question to be $\frac{1}{2} + \frac{1}{2} \ln\left(\frac{1}{2}\right)$.

$$\begin{aligned}
 \int_0^{1/2} \int_0^{1/2} \frac{1}{1-x} dy dx &= \frac{1}{2} \int_0^{1/2} \frac{1}{1-x} dx \\
 &= \frac{1}{2} [-\ln(1-x)]_0^{1/2} \\
 &= -\frac{1}{2} \ln\left(\frac{1}{2}\right)
 \end{aligned}$$

Subtracting the lower triangle from the square, we see that the integral over the horizontally hatched triangle equals

$$-\frac{1}{2} \ln\left(\frac{1}{2}\right) - \frac{1}{2} - \frac{1}{2} \ln\left(\frac{1}{2}\right) = -\frac{1}{2} - \ln\left(\frac{1}{2}\right)$$

We can then calculate

$$P\left(X < \frac{1}{2}, Y < \frac{1}{2} \mid X + Y > \frac{1}{2}\right) = \frac{-\frac{1}{2} - \ln\left(\frac{1}{2}\right)}{\frac{1}{2} - \frac{1}{2} \ln\left(\frac{1}{2}\right)} \approx 0.23$$

s.7.3.6. a. Since $f_{XY}(x, y)$ is a joint probability density function, we should have

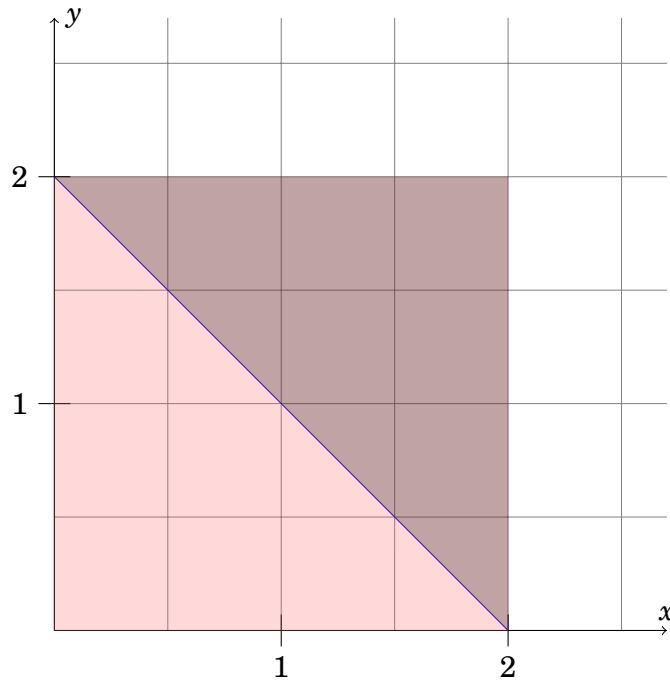
1. $f_{X,Y}(x, y) \geq 0$. This is satisfied since $x, y \geq 0$.

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{XY}(x, y) dx dy = 1,$$

. We can calculate the integral as follows.

$$\begin{aligned}
 \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f_{XY}(x, y) dx dy = 1 &\iff \frac{3}{16} \int_0^c \int_0^c xy^2 dx dy = 1 \\
 &\iff \frac{3}{32} \int_0^c (x^2 y^2 \Big|_0^c) dy = 1 \\
 &\iff \frac{3c^2}{32} \int_0^c y^2 dy = 1 \\
 &\iff \frac{3c^2}{96} y^3 \Big|_0^c = 1 \\
 &\iff 3c^5 = 96 \iff c = 2
 \end{aligned}$$

b. First, draw the area over which the integral is taken.



We want to integrate over the darkest area. Hence, for every value of y , x varies between $2 - y$ and 2. Hence, the required probability can be calculated as follows:

Solution 1

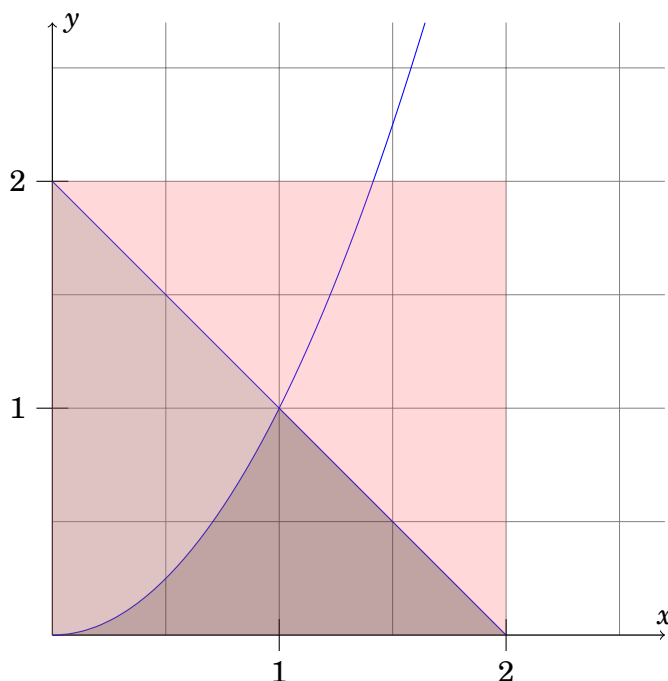
$$\begin{aligned}
P(X+Y > 2) &= \\
&= \int_0^2 \int_{2-y}^2 f_{X,Y}(x,y) dx dy \\
&= \frac{3}{16} \int_0^2 \int_{2-y}^2 xy^2 dx dy \\
&= \frac{3}{32} \int_0^2 (x^2 y^2 |_{2-y}^2) dy \\
&= \frac{3}{32} \int_0^2 (4y^2 - (2-y)^2 y^2) dy \\
&= \frac{3}{32} \int_0^2 (4y^3 - y^4) dy \\
&= \frac{3}{32} \left(y^4 - \frac{1}{5} y^5 \right) \Big|_0^2 \\
&= \frac{3}{32} \left(16 - \frac{32}{5} \right) \\
&= \frac{3}{2} - \frac{3}{5} = \frac{9}{10}
\end{aligned}$$

draw the area over which the integral is taken.

Solution 2

$$\begin{aligned}
P(X+Y > 2) &= \\
&= \int_0^2 \int_{2-x}^2 f_{X,Y}(x,y) dy dx \\
&= \frac{3}{16} \int_0^2 \int_{2-x}^2 xy^2 dy dx \\
&= \frac{3}{16} \int_0^2 \left(\frac{1}{3} xy^3 \Big|_{2-x}^2 \right) dx \\
&= \frac{1}{16} \int_0^2 (8x - x(2-x)^3) dx \\
&= \frac{1}{16} \int_0^2 (x^4 + 12x^2 - 6x^3) dx \\
&= \frac{1}{16} \left(\frac{1}{5} x^5 + 4x^3 - \frac{3}{2} x^4 \right) \Big|_0^2 \\
&= \frac{1}{16} \left(\frac{32}{5} + 32 - 24 \right) \\
&= \frac{2}{5} + \frac{8}{16} = \frac{9}{10}
\end{aligned}$$

c. First,



The conditional probability is given by

$$P(Y < X^2 | X+Y < 2) = \frac{P(X+Y < 2 \cap Y < X^2)}{P(X+Y < 2)}.$$

From part (b) we have $P(X + Y < 2) = 1 - P(X + Y > 2) = \frac{1}{10}$, which means the probability of falling into one of the two darkest areas equals $\frac{1}{10}$. The probability $P(X + Y < 2 \cap Y < X^2)$ is given by the integral over the darkest area in the plot.

Solution 1

$$\begin{aligned}
 P(X + Y < 2 \cap Y < X^2) &= \int_0^1 \int_0^{x^2} f_{X,Y}(x,y) dy dx + \int_1^2 \int_0^{2-x} f_{X,Y}(x,y) dy dx \\
 &= \frac{3}{16} \left[\int_0^1 \int_0^{x^2} xy^2 dy dx + \int_1^2 \int_0^{2-x} xy^2 dy dx \right] \\
 &= \frac{1}{16} \left[\int_0^1 \left(xy^3 \Big|_0^{x^2} \right) dx + \int_1^2 \left(xy^3 \Big|_0^{2-x} \right) dx \right] \\
 &= \frac{1}{16} \left[\int_0^1 x^7 dx + \int_1^2 x(2-x)^3 dx \right] \\
 &= \frac{1}{16} \left[\int_0^1 x^7 dx + \int_1^2 (-x^4 + 6x^3 - 12x^2 + 8x) dx \right] \\
 &= \frac{1}{128} \left[x^8 \Big|_0^1 \right] + \frac{1}{16} \left(-\frac{1}{5}x^5 + \frac{3}{2}x^4 - 4x^3 + 4x^2 \right) \Big|_1^2 \\
 &= \frac{1}{128} + \frac{1}{16} \left(-\frac{32}{5} + 24 - 32 + 16 + \frac{1}{5} - \frac{3}{2} + 4 - 4 \right) = \frac{17}{640}
 \end{aligned}$$

Solution 2 (easier)

$$\begin{aligned}
 P(X + Y < 2 \cap Y < X^2) &= \int_0^1 \int_{\sqrt{y}}^{2-y} f_{X,Y}(x,y) dx dy \\
 &= \frac{3}{16} \left[\int_0^1 \int_{\sqrt{y}}^{2-y} xy^2 dx dy \right] \\
 &= \frac{3}{32} \left[\int_0^1 \left(x^2 y^2 \Big|_{\sqrt{y}}^{2-y} \right) dy \right] \\
 &= \frac{3}{32} \left[\int_0^1 (2-y)^2 y^2 - y^3 dy \right] \\
 &= \frac{3}{32} \left[\int_0^1 4y^2 - 5y^3 + y^4 dy \right] \\
 &= \frac{3}{32} \left[\frac{4}{3}y^3 - \frac{5}{4}y^4 + \frac{1}{5}y^5 \Big|_0^1 \right] \\
 &= \frac{3}{32} \left(\frac{4}{3} - \frac{5}{4} + \frac{1}{5} \right) = \frac{17}{640}
 \end{aligned}$$

Using either Solution 1 or Solution 2, we get the final answer:

$$P(Y < X^2 | X + Y < 2) = \frac{P(X + Y < 2 \cap Y < X^2)}{P(X + Y < 2)} = \frac{\frac{17}{640}}{\frac{1}{10}} = \frac{17}{64}.$$

s.7.4.2. Use the hint. For $k = 0, 1, \dots, n$,

$$P\{X = k, N = n\} = P\{X = k, Y = n - k\} = pq^k pq^{n-k} = p^2 q^n.$$

(Have we used independence somewhere?) Now observe that the right hand side does not depend on k . This implies that $P\{X = k|N = n\}$ also does not depend on k . (Why?) But, since $P\{X = k|N = n\}$ is a true PMF, it must be that $\sum_{k=0}^n P\{X = k|N = n\}$ adds up to 1. These two ideas put together imply that $P\{X = k|N = n\} = 1/(n + 1)$.

With Bayes' expression

$$P\{X = k|N = n\} = \frac{P\{X = k, N = n\}}{P\{N = n\}},$$

it follows that

$$P\{N = n\} = \frac{P\{X = k, N = n\}}{P\{X = k|N = n\} P\{N = n\}} = (n + 1)p^2 q^n.$$

s.7.4.3. Just reasoning as if there is no problem, i.e., applying Bayes' rule in a naive way,

$$\begin{aligned} F_T(t|x) &= P\{T \leq t|X = x\} = P\{X + Y \leq t|X = x\} \\ &= P\{Y \leq t - x, X = x\} / P\{X = x\} = P\{Y \leq t - x\} P\{X = x\} / P\{X = x\} \\ &= P\{Y \leq t - x\}, 0 \leq x \leq t. \end{aligned}$$

where I use that Y and X are independent to split the probability.

The problem with this derivation is that we multiply and divide by 0 ($= P\{X = x\}$) just as if all is ok. But hopefully, you know that when we multiply and divide by zero, we can get any answer we like. A better way is as follows. Note beforehand that I do not expect that you could have come up with such an answer, but you should definitely study it.

The first step is to realize that PDF $f_{T|X}(t|x) = f_{TX}(t, x)/f_X(x)$ is well defined; we don't divide by zero because $f_X(x) > 0$ on $x \geq 0$. By the proof of BH.8.2.1 we see that $f_{TX}(t, x) = f_X(x)f_Y(t - x)I_{0 \leq x \leq t}$, where I include the indicator to ensure that we don't run out of the support of X and T . Thus,

$$f_{T|X}(t|x) = \frac{f_{TX}(t, x)}{f_X(x)} = f_X(x)f_Y(t - x)I_{0 \leq x \leq t}/f_X(x) = f_Y(t - x)I_{0 \leq x \leq t}.$$

Now we know that a conditional PDF is a full-fledged PDF. So we can use idea that to *define* the conditional CDF as follows:

$$F_{T|X}(t|x) := \int_0^t f_{T|X}(v|x)I_{0 \leq x \leq v} dv = \int_x^t \lambda e^{-\lambda(v-x)} dv = \int_0^{t-x} \lambda e^{-\lambda v} dv = 1 - e^{-\lambda(t-x)}.$$

Isn't it a bit strange that we get the same answer? How to get out of this situation in a technically correct way is one of the hard parts of (mathematical) probability, and certainly not something we can deal with in this course. All books on elementary¹ probability, and lecturers similarly, struggle with this problem; this course is not an exception, nor am I.

¹ When in mathematics something is elementary, it doesn't necessarily mean that that thing is simple. In fact, it can be very difficult. Elementary means that we just don't use very advanced mathematical concepts.)

- b. See part a.
c. By the above,

$$\begin{aligned} f_{X|T}(x|t) &= f_{TX}(t, x)/f_T(t), \\ f_{TX}(t, x) &= f_X(x)f_Y(t-x)I_{0 \leq x \leq t} \lambda^2 e^{-\lambda x} e^{-\lambda(t-x)} I_{0 \leq x \leq t} = \lambda^2 e^{-\lambda t} I_{0 \leq x \leq t}, \\ \implies f_{X|T}(x|t) &\propto \lambda^2 e^{-\lambda t} I_{0 \leq x \leq t}, \end{aligned}$$

where the last follows because $f_T(t)$ is just a normalization constant. Now we use some real nice, but subtle, reasoning to avoid computing f_T by means of marginalizing out x from $f_{T,X}(t, x)$. Observe that $f_{TX}(t, x)$ is constant *as a function of x* on $0 \leq x \leq t$ (in other words, the RHS does not depend on x on this interval). But $f_{X|T}(x|t)$ is also a real PDF. This implies that the constant $f_T(t)$ (since it does not depend on x) must be such that $f_{X|T}(x|t)$ integrates to 1 on $0 \leq x \leq t$. The only possibility is that $f_{X|T}(x|t) = t^{-1} I_{0 \leq x \leq t}$.

This reasoning gives some offspin. We can conclude that

$$f_T(t) = f_{TX}(t, x)/f_{T|X}(t, x) = \lambda^2 t e^{-\lambda t}.$$

This is more than a nice trick. Recall it, as it is not only used more often in the book, but also in more advanced courses on data science and machine learning.

s.7.4.4. c. Here is the answer. The ideas are important, you'll need them during nearly any course in statistics, given the importance of the normal distribution.

$$f_{Y,Z}(y, z) = \int \frac{1}{2\pi} e^{-(y-x)^2/2} e^{-(z-x)^2/2} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx.$$

It remains to simplify $(y-x)^2 + (z-x)^2 + x^2$. With a bit of work, it follows that this can be written as

$$3(x - (y+z)/3)^2 - (y+z)^2/3 + y^2 + z^2.$$

When plugging this in the integral, the last two terms appear in front of the integral. The term $(y+z)/2$ is just a shift, hence can be neglected in the integration over x . The 3 has to be absorbed in the standard deviation $\sigma = 1/\sqrt{3}$. And therefore,

$$f_{Y,Z}(y, z) = \frac{1}{2\pi} \frac{1}{\sqrt{3}} e^{-y^2/2 - z^2/2 + (y+z)^2/6}.$$

s.7.4.5.

s.7.4.6. Using the hint: $F(x, y)$ is the area of an (infinite) square lying south west of the point (x, y) . Add and subtract such (infinite) squares until the square $[a_1, a_2] \times [b_1, b_2]$ is covered exactly once. Realize that in the process, the square $(-\infty, a_1] \times (-\infty, b_1]$ is subtracted twice.

s.7.4.7. From the hint,

$$\begin{aligned} P\{Y_1 < cY_2\} &= \int \int I_{x < cy} \lambda_1 e^{-\lambda_1 x} \lambda_2 e^{-\lambda_2 y} dx dy = \lambda_1 \lambda_2 \int_0^\infty e^{-\lambda_1 x} \int_{x/c}^\infty e^{-\lambda_2 y} dy dx \\ &= \lambda_1 \int_0^\infty e^{-\lambda_1 x} e^{-\lambda_2 x/c} dx = \frac{\lambda_1}{\lambda_1 + \lambda_2/c}. \end{aligned}$$

Check the result for $c = 0$ and $c = \infty$.

I prefer to use conditioning, like this:

$$\begin{aligned} P\{Y_1 < cY_2\} &= \int P\{Y_1 < cY_2 | Y_1 = x\} \lambda_1 e^{-\lambda_1 x} dx = \int P\{Y_2 > x/c | Y_1 = x\} \lambda_1 e^{-\lambda_1 x} dx \\ &= \int e^{-\lambda x/c} \lambda_1 e^{-\lambda_1 x} dx, \end{aligned}$$

and the rest goes as before. Actually, I tend to use conditioning as it helps to make the reasoning easier. In this case, suppose that I know that $Y_1 = x$, what can I say about $P\{Y_2 > cx\}$?

BTW, conditioning does not always make things simpler. When rvs are dependent, then you have to watch out.

s.7.4.8.

s.7.4.9. First check [7.1.28].

In general, I am always very careful with such ‘shortcuts’ such as $\max\{X, Y\} + \min\{X, Y\} = X + Y$. As a matter of fact, I try to avoid such arguments because it is easy to go wrong. Seemingly plausible arguments are often wrong due to overlooked dependency or non-linearity (effects of higher moments).

It is useful to write $\max\{x, y\} = xI_{x \geq y} + yI_{y > x}$, and something similar for the minimum. In the present case, $\text{Cov}[X, Y] = E[XY] - E[X]E[Y]$, and, similarly, $\text{Cov}[M, L] = E[ML] - E[M]E[L]$, where M is max, and L is min. With the above indicators, it is simple to show that $E[ML] = E[XY]$:

$$\begin{aligned} ML &= (XI_{X \geq Y} + YI_{Y \geq X})(XI_{X < Y} + YI_{Y < X}) \\ &= XYI_{X \geq Y} + XYI_{Y < X} = XY \end{aligned}$$

since $I_{X \geq Y}I_{X < Y} = 0$.

However, take $X, Y \sim \text{Exp}(\lambda)$. Then, $E[M] = 3/(2\lambda)$ and $E[L] = 1/(2\lambda)$, but $E[X] = E[Y] = 1/\lambda$.

s.7.4.10. a. In my notation, $X_i = 0 \implies Y_i \neq 0$ and $X_i \neq 0 \implies Y_i = 0$. The reason is that in step i , the drunkard makes a step left or right OR up or down. However, s/he cannot move to the right and up at the same time.

Here is an argument based on recursion. (By now I hope you see that I like this method in particular).

$$E[R_n^2] = E[(R_{n-1} + X_n + Y_n)^2],$$

but R_{n-1} and $X_n + Y_n$ are independent, and $E[(X_n + Y_n)^2] = 1$. Using the recursion, $E[R_n^2] = n$.

s.7.4.11. b. Use the hint. Then, if we choose two points at random from the sample, then $(x_i - x_j)(y_i - y_j)$ is the area spanned by these two points. More generally, I have n choices for my first point, and also n choices for the second point (if both points are the same, the area of the rectangle is 0, so we don't have to exclude such choices). Hence, the expected area of the rectangle spanned by the two random points (X, Y) and (\tilde{X}, \tilde{Y}) is

$$\frac{1}{n^2} \sum_{i,j} (x_i - x_j)(y_i - y_j).$$

Simplify this to show that

$$2 \frac{1}{n} \sum_i x_i y_i - 2\bar{x}\bar{y} = 2r$$

Hence, by part a., the expected area is twice the covariance.

Why is $\text{Cov}[X, a] = 0$ for a a constant? Because the 'area' of rectangles, all with the same y -coordinate, is zero, i.e., they lie on a line.

c. This is the part of the exercise that explains what the above is all about. Since there is a direct relation between covariance and area, we can use geometric arguments to derive (and memorize!) all properties of covariance! Write property i. of covariance as $\text{Cov}[X, Y] = \text{Cov}[Y, X]$. Suppose I flip the x and y -axis, does the area of a rectangle change? For property ii., what happens to the area of rectangle if you stretch the sides? For property iii., realize that this is just a shift of a rectangle that leaves its area invariant. For property iv., what happens to the area if you put an extra rectangle on top or to the right?

BTW, property iii. follows directly from property iv. In iv., take W_3 equal to a constant a_2 , in other words $P\{W_3 = a_2\} = 1$. We know that $\text{Cov}[X, a] = 0$ for a constant a .

Here are some final remarks.

Let's put all the above in a very general frame. The covariance has a number of interesting properties:

1. It is bilinear, that is, the covariance is linear in both arguments. The linearity in the first argument means that $\text{Cov}[X + Y, Z] = \text{Cov}[X, Z] + \text{Cov}[Y, Z]$ and $\text{Cov}[aX, Z] = a \text{Cov}[X, Z]$ for $a \in \mathbb{R}$. The linearity in the second argument means that $\text{Cov}[X, Y + Z] = \text{Cov}[X, Y] + \text{Cov}[X, Z]$ and $\text{Cov}[X, aZ] = a \text{Cov}[X, Z]$ for $a \in \mathbb{R}$.
2. It is symmetric: $\text{Cov}[X, Y] = \text{Cov}[Y, X]$, from which we define $V[X] = \text{Cov}[X, X]$.
3. $\text{Cov}[X, a] = 0$ for all $a \in \mathbb{R}$.

If you memorize the first two properties of covariance, all the rest follows.

Now we do some geometry. Take three vectors $x, y, z \in \mathbb{R}^2$ (it's easy to generalize to \mathbb{R}^n). Then we know that the area $D(x, y)$ of the parallelogram spanned by vectors x and y satisfies the following properties.

1. Area is bilinear. The linearity in the first argument means that $D(x + y, z) = D(x, z) + D(y, z)$ and $D(ax, z) = aD(x, z)$ for $a \in \mathbb{R}$. (Just make a drawing to convince you about this.) The linearity in the second argument means that $D(x, y + z) = D(x, y) + D(x, z)$ and $D(x, az) = aD(x, z)$ for $a \in \mathbb{R}$.

2. $D(x, x) = 0$; there is no area between x and x .
3. $D((1, 0), (0, 1)) = 1$; the area of the square with side 1 is 1.

In fact, the first property means that stretching vectors and stacking parallelograms result in stretching and adding areas. The second says that the area of a parallelogram spanned by two parallel vectors is zero. The third specifies that the area of the unit square is 1.

Now it can be proven that there exists just one function D that satisfies these properties. In fact, this is the determinant of the matrix with as columns the vectors that span the parallelogram. Moreover, it can be shown that the second property can be replaced by the skew-symmetric property: $D(x, y) = -D(y, x)$. (Note that $D(x, x) = -D(x, x) \implies 2D(x, x) = 0 \implies D(x, x) = 0$.)

Let us use the properties to compute the area of a parallelogram spanned by the vectors $x = (a, b)$ and $y = (c, d)$ in 2D. Then

$$\begin{aligned} D(x, y) &= D((a, b), (c, d)) = D(a(1, 0) + b(0, 1), c(1, 0) + d(0, 1)) \\ &= adD((1, 0), (0, 1)) + bcD((0, 1), (1, 0)) = ad - bc, \end{aligned}$$

where we use bilinearity in the first step, and skew-symmetry in the second and third. And this is indeed the determinant of the matrix with x and y as columns.

So, all in all, this is what I remembered throughout the years: the covariance and the determinant are bi-linear forms, the first is symmetric, the second skew- (or anti-)symmetric.

Finally, I don't see why the areas of the rectangles have to have a sign in this problem. Interestingly, for the determinant, the areas of the parallelograms do have to have a sign to make the concept useful for physics.

s.7.4.12. a. Follows directly from the hint.

Check the hint!

c. If X_i and Y_j are iid, it must be that $w_1 = n/(n + m)$.

b. Can we make some further progress, just by keeping a clear mind? Well, in fact we can by using our insights of part c. If we have $n + m$ iid measurements of which we call n measurements of type X_i , and m of type Y_j , then

$$V[\hat{\theta}_1] = E\left[\left(\frac{1}{n} \sum_i X_i - \theta\right)^2\right] = n^{-2} E\left[\left(\sum_i (X_i - \theta)\right)^2\right] = n^{-2} V\left[\sum_i X_i\right] = V[X_1]/n = \sigma^2/n.$$

So, $n = \sigma^2/V[\hat{\theta}_1]$, and likewise $m = \sigma^2/V[\hat{\theta}_2]$. Finally, plug this into our earlier expression for w_1 to get

$$w_1 = \frac{n}{n + m} = \frac{\sigma^2/V[\hat{\theta}_1]}{\sigma^2/V[\hat{\theta}_1] + \sigma^2/V[\hat{\theta}_2]} = \frac{V[\hat{\theta}_2]}{V[\hat{\theta}_1] + V[\hat{\theta}_2]}.$$

If we check our earlier insight, then we see that if $V[Y_j] = 0$, then $V[\hat{\theta}_2] = 0$, hence $w_1 = 0$ in that case. This is precisely what we wanted.

Let us finally use the hint of BH to check that the above expression for w_1 is correct.

$$E[(\hat{\theta} - \theta)^2] = E[(w_1(\hat{\theta}_1 - \theta) + w_2(\hat{\theta}_2 - \theta))^2] = V[w_1\hat{\theta}_1] + V[w_2\hat{\theta}_2],$$

by independence. Take the w 's out of the variances, then write $w_2 = 1 - w_1$, take ∂_{w_1} of the expression, set the result to 0, and solve for w_1 . You'll get the above expression.

s.7.4.13. a. Multinomial.

b. With the hint we end up at $X_1 + X_2 \sim \text{Bin}(n, p^2 + 2p(1 - p))$.

c. Here is a short intermezzo on finding a recursion for the sum of a number of Bernoulli rvs. Let S_n be the number of successes in the binomial, and write $g_n(i) = P\{S_n = i\}$ for this case. Then,

$$\begin{aligned} g_n(i) &= g_{n-1}(i-1)p + g_{n-1}(i)q \\ &= (g_{n-2}(i-2)p + g_{n-2}(i-1)q)p + (g_{n-2}(i-1)p + g_{n-2}(i)q)q \\ &= g_{n-2}(i-2)p^2 + g_{n-2}(i-1)2pq + g_{n-2}(i)q^2. \end{aligned}$$

I also know that $g_n(i) = \binom{n}{i} p^i q^{n-i}$. End of intermezzo.

Now compare the recursion with $f_n(i)$ for the genes to the expression for the binomial. They are nearly the same, except that in the genes case, the 'n' seems to run twice as fast. I then tried the guess $f_n(i) = \binom{2n}{i} p^i q^{2n-i}$. For you, plug it in, and show that it works.

So, what was my overall approach? I used recursion, but got stuck. Then I used recursion for a simpler case whose solution I know by heart. I compared the recursions for both cases to see whether I could recognize a pattern. This led me to a guess, which I verified by plugging it in. Using recursion is not guaranteed to work, of course, but often it's worth a try.

Now, looking back, I realize that it is as if individual n adds the outcome of two coin flips (with values in AA , Aa or aa) to the sum S_n of A 's. For you to solve: what is the distribution of two coin flips? Next, S_n is just the sum of n individual 'double coin flips'. Hence, what must the distribution of S_n be?

d. It is easiest to work with $f(p) = \log P\{X_1 = k, X_2 = l, X_3 = m\}$. With part a. this can be written as

$$f(p) = C + (2k + l)\log p + (l + 2m)\log(1 - p),$$

where C is a constant (the log of the normalization constant). (BTW, with this you can check your answer for part a.) Compute $df(p)/dp = 0$, because at this p , $\log f$, hence f itself, is maximal. Observe that C drops out of the computation, because when differentiating, it disappears.

e. Now we like to know what p maximizes $P\{X_3 = n - i\}$. Take $g(q) = \log P\{X_3 = n - i\}$, then

$$g(q) = C + i\log(1 - q^2) + 2(n - i)\log q.$$

(With this, check your answer of part b.) Again, take the derivative (with respect to q), and solve for q .

s.7.4.14. a. It is given that $P\{T \leq t | D = 1\} = G(t)$ and $P\{T \leq t | D = 0\} = H(t)$. From Theorem 5.3.1.i, we have that we can associate a rv. to a CDF F . Sometimes we say that the CDF F /induces/ a rv. X . So let us use this here to say that G induces the rv. T_1 and H induces T_0 . So the /sensitivity/ is $P\{T_1 > t_0\} = 1 - G(t_0)$ and the /specificity/ is $P\{T_1 < t_0\} = H(t_0)$.

To make the ROC plot, I first made two plots, one of the sensitivity and the other for 1 minus the specificity, i.e., $1 - H(t_0)$. Then, in the ROC plot, we put a specificity of s on the x -axis, then we search for a t such that $1 - H(t) = s$, and then we plug this t into $1 - G(t)$ to get the sensitivity. To help you understand this better, check that $s = 0 \implies t = b \implies 1 - G(t) = 0$. Moreover, check that $s = 1 \implies t = a \implies 1 - G(t) = 1$. Hence, the ROC curve starts in the origin and stops at the point $(1, 1)$.

With this insight, the area under the ROC curve can be written as

$$\int_0^1 (1 - G(H^{-1}(1 - s))) ds = 1 - \int_0^1 G(H^{-1}(1 - s)) ds = 1 - \int_a^b G(t)h(t) dt,$$

where, in the last step, we use the 1D change of variable $H(t) = 1 - s \implies h(t)dt = -ds$. It remains to interpret the integral, so let's plug in the definitions:

$$\int_a^b G(t)h(t) dt = \int_a^b P\{T_1 \leq t\} f_{T_0}(t) dt = \int_a^b P\{T_1 \leq T_0 | T_0 = t\} f_{T_0}(t) dt = P\{T_1 \leq T_0\}.$$

s.8.1.1. Recall that $V[X] = E[X^2] - (E[X])^2$; so we have to deal with the function $g(x) = x^2$ because $E[X^2] = E[g(X)]$. Note that even to properly define the variance, we have to deal with a function that is not one-to-one everywhere on \mathbb{R} .

s.8.1.2.

$$X \in \{0, \dots, 5\} \implies Z \in \{0, 3, 6, 9, 12, 15\}, \quad \text{and not in } \{0, 1, 2, \dots, 14, 15\}, \quad (13.0.59)$$

$$z = g(x) = 3x, \quad (13.0.60)$$

$$p_Z(z) = \sum_{x: g(x)=z} p_X(x) = \frac{1}{6} I_{z \in \{0, 3, 6, 9, 12, 15\}}, \quad (13.0.61)$$

$$F_Z(z) = \frac{1}{6} \sum_{x=0}^z I_{x \in \{0, 3, 6, 9, 12, 15\}}. \quad (13.0.62)$$

s.8.1.3. To get the derivative of g^{-1} , consider the equality $g(g^{-1}(y)) = y$. Then, taking derivatives with respect to y at both sides, and applying the chain rule,

$$g(g^{-1}(y)) = y \implies \frac{d}{dy} g(g^{-1}(y)) = 1 \iff g'(g^{-1}(y)) \frac{d}{dy} g^{-1}(y) = 1 \implies \frac{d}{dy} g^{-1}(y) = 1/g'(x),$$

where we use that $g^{-1}(y) = x$.

s.8.1.4. When working the CDFs, we need to solve the problem $\{x : g(x) \leq y\}$. If we take $g(x) = \sin x + x/100$ then this is really messy. In fact, to solve this, we first solve for the set $x : g(x) = y$, which might still be hard, but requires less work than check each and every interval.

With PDFs we only have to require *locally* that g is one-to-one, and we don't have to work with inequalities, but can directly focus on the set $\{x : g(x) = y\}$.

In 2D, functions can have saddle points, i.e., points in which the function increases in one direction and decreases in another. Then finding the set of points x such that $g(x, y) \leq (u, v)$ (which we need if we want to express $P\{g(X, Y) \leq (u, v)\}$ in terms of the distribution $F_{X,Y}$) is not a particularly attractive task, to say the least.

s.8.1.5. Note that $g(x) = x^2$ is not monotone increasing, moreover, $g^{-1}(y)$ does not exist (in \mathbb{R}) for $y < 0$. We split the line into disjoint intervals in which g is either strictly increasing or decreasing, and then we apply the above rule in each of the intervals. Since $g'(x) = 2x$ and $x = \pm\sqrt{y}$,

$$f_Y(y) = f_X(\sqrt{y}) \frac{1}{2\sqrt{y}} + f_X(-\sqrt{y}) \frac{1}{2\sqrt{y}}.$$

s.8.1.6. Take $y = g(x) = \lambda x$. Then,

$$f_Y(y) = f_X(x) \frac{dx}{dy} = f_X(x) \frac{1}{g'(x)} = e^{-y/\lambda} \frac{1}{\lambda}.$$

With this,

$$E[Y] = \int_0^\infty y f_Y(y) dy = \int_0^\infty y e^{-y/\lambda} \frac{1}{\lambda} dy.$$

To solve this integral, I recognize y/λ in the exponent, and I want to get rid of the $1/\lambda$ factor. Hence, I write $u = y/\lambda$, and use this to see that

$$u = y/\lambda \implies du = dy/\lambda \implies dy = \lambda du.$$

Then, including a and b for the boundaries to show explicitly what is going on when changing the variables

$$\int_a^b y/\lambda e^{-y/\lambda} dy = \int_{a/\lambda}^{b/\lambda} u e^{-u} \lambda du = \lambda \int_{a/\lambda}^{b/\lambda} u e^{-u} du.$$

Applying this to our case so that $a = 0/\lambda = 0$ and $b = \infty/\lambda = \infty$,

$$E[Y] = \lambda \int_0^\infty u e^{-u} du = \lambda E[X].$$

s.8.1.7. When we have the density f_Y and the function g , then the substitution rule says that,

$$\int_a^b f_Y(g(x)) g'(x) dx = \int_{g(a)}^{g(b)} f_Y(y) dy.$$

We also want that the transformation from X to Y does not affect the probability of the set (event) $A = [a, b]$, hence,

$$\int_{g(a)}^{g(b)} f_Y(y) dy = \int_a^b f_X(x) dx.$$

Combining the above two equations gives that

$$\int_a^b f_Y(g(x))g'(x) dx = \int_a^b f_X(x) dx.$$

Since this holds for any a and b , it follows that

$$f_Y(g(x))g'(x) = f_X(x).$$

s.8.1.8.

s.8.1.9.

s.8.1.10.

s.8.1.11.

$$X \in [0, 5] \implies Z \in [0, 15], \quad (13.0.63)$$

$$z = 3x = g(x) \implies x = z/3, \quad (13.0.64)$$

$$f_Z(z) = f_X(x) \frac{dx}{dz}, \quad (13.0.65)$$

$$\frac{dz}{dx} = 3, \quad (13.0.66)$$

$$f_Z(z) = f_X(z/3) \frac{1}{3}. \quad (13.0.67)$$

$F_Z(u) = 1$ for $u \geq 15$ and $F_Z(u) = 0$ for $u \leq 0$. When $0 \leq u \leq 15$,

$$F_Z(u) = \int_0^u f_X(z/3) \frac{1}{3} dz = \frac{1}{5} \int_0^u I_{0 \leq z/3 \leq 5} \frac{1}{3} dz \quad (13.0.68)$$

$$= \frac{1}{5} \int_0^u I_{0 \leq z \leq 15} \frac{1}{3} dz = \frac{u}{15}. \quad (13.0.69)$$

s.8.1.12.

$$X \in [0, 5] \implies Z \in [0, 125], \quad (13.0.70)$$

$$z = x^3 = g(x) \implies x = z^{1/3}, \quad (13.0.71)$$

$$f_Z(z) = f_X(x) \frac{dx}{dz}, \quad (13.0.72)$$

$$\frac{dz}{dx} = 3x^2 = 3z^{2/3}, \quad (13.0.73)$$

$$f_Z(z) = f_X(z^{1/3}) \frac{1}{3z^{2/3}}. \quad (13.0.74)$$

When $F_Z(u) = 1$ for $u \leq 125$ and $F_Z(u) = 0$ for $u \leq 0$. When $0 \leq u \leq 125$,

$$F_Z(u) = \int_0^u f_X(z^{1/3}) \frac{1}{3z^{2/3}} dz = \frac{1}{5} \int_0^u I_{0 \leq z^{1/3} \leq 5} \frac{1}{3z^{2/3}} dz \quad (13.0.75)$$

$$= \frac{1}{5} \int_0^u I_{0 \leq z \leq 125} \frac{1}{3z^{2/3}} dz \quad (13.0.76)$$

$$= \frac{1}{5} \int_0^u \frac{1}{3z^{2/3}} dz = \frac{1}{5} z^{1/3} \Big|_0^u = u^{1/3}/5. \quad (13.0.77)$$

s.8.1.13.

$$z = g(x) = (x - \mu)/\sigma, \implies x = \sigma z + \mu \quad (13.0.78)$$

$$f_Z(z) = f_X(x) \frac{dx}{dz}, \quad (13.0.79)$$

$$\frac{dz}{dx} = \frac{1}{\sigma}, \quad (13.0.80)$$

$$f_Z(z) = f_X(x)\sigma = \sigma f_X(\sigma z + \mu) \quad (13.0.81)$$

and now using the density of $X \sim \text{Norm}(\mu, \sigma)$,

$$f_Z(z) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(\sigma z + \mu - \mu)^2/2\sigma^2} \sigma = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}. \quad (13.0.82)$$

s.8.1.14.

$$z = g(x) = e^{-x} \implies x = -\log z, \quad (13.0.83)$$

$$x \in (0, \infty) \implies z \in (0, 1), \quad (13.0.84)$$

$$f_Z(z) = f_X(x) \frac{dx}{dz}, \quad (13.0.85)$$

$$\frac{dz}{dx} = -e^{-x}, \quad \text{Don't forget to take the abs value next,} \quad (13.0.86)$$

$$f_Z(z) = f_X(x)e^x = e^{-x}e^x = 1 I_{0 < z < 1}, \quad (13.0.87)$$

where we include the domain of Z in the last equality.

s.8.1.15.

$$(u, v) = (x + y, x - y) = g(x, y) \implies (x, y) = ((u + v)/2, (u - v)/2), \quad (13.0.88)$$

$$\frac{\partial(u, v)}{\partial(x, y)} = \begin{vmatrix} 1 & 1 \\ 1 & -1 \end{vmatrix} = -2 \implies |-2| = 2, \quad (13.0.89)$$

$$f_{U,V}(u, v) = f_{X,Y}(x, y) \frac{\partial(x, y)}{\partial(u, v)} = f_{X,Y}((u + v)/2, (u - v)/2)/2 \quad (13.0.90)$$

$$= \frac{1}{4\pi} e^{-((u+v)/2)^2/2} e^{-((u-v)/2)^2/2} \quad (13.0.91)$$

$$= \frac{1}{4\pi} e^{-u^2/4 - v^2/4}, \quad (13.0.92)$$

where we work out the squares and simplify. Hence, U and V are independent and normally distributed with mean 0 and $\sigma = \sqrt{2}$. This is in line with our earlier definition of a multi-variate normal distribution.

- s.8.1.16.** 1. $Z = Y^4 \in [0, \infty)$ for $Y \in (-\infty, \infty)$;
 2. $Y = X^3 + a \in (a, a + 1)$ for $X \in (0, 1)$;
 3. $U = |V| + b \in [b, \infty)$ for $V \in (-\infty, \infty)$;
 4. $Y = e^{X^3} \in (0, \infty)$ for $X \in (-\infty, \infty)$;
 5. $V = U I_{U \leq c} \in (-\infty, c]$ for $U \in (-\infty, \infty)$;
 6. $Y = \sin(X) \in [-1, 1]$ for $X \in (-\infty, \infty)$;
 7. $Y = \frac{X_1}{X_1 + X_2} \in (0, 1)$ for $X_1 \in (0, \infty)$ and $X_2 \in (0, \infty)$;
 8. $Z = \log(UV) \in (-\infty, \infty)$ for $U \in (0, \infty)$ and $V \in (0, \infty)$.

s.8.1.17. When the variables become dependent, the Jacobian becomes zero. For instance, in the latter case,

$$\frac{\partial(u, v)}{\partial(x, y)} = \begin{vmatrix} 1/y & -x/y^2 \\ -y/x^2 & 1/x \end{vmatrix} = \frac{1}{xy} - \frac{x}{y^2} \frac{y}{x^2} = 0. \quad (13.0.93)$$

Moreover, the function g is not locally one-to-one.

s.8.1.18. If we would not add this extra variable, we cannot use the change of variables theorem. We also need a function to deal with the scaling. In the change of variables theorem, this is the Jacobian.

There is also another problem. Consider the function $g(x, y)$ that maps \mathbb{R}^2 to \mathbb{R} . The inverse set $\{(x, y) : g(x, y) = z\}$ can be quite complicated, while the set $\{y : g(x, y) = z\}$ for a fixed x is hopefully just one point. Hence, the mapping $(x, y) \rightarrow (x, g(x, y))$ is, at least locally, one-to-one.

It is possible to deal with the more general problem, but this requires much more theory than we need for this course.

s.8.1.19. From BH.8.1.4: Z chi-square $\implies X = \sqrt{Z} \sim \text{Norm}(0, 1)$. Then, from BH.8.1.9,

$$X^2 + Y^2 = (\sqrt{2T} \cos U)^2 + (\sqrt{2T} \sin U)^2 = 2T (\cos^2 U + \sin^2 U) = 2T \sim \text{Exp}(1/2), \quad (13.0.94)$$

when $X, Y \sim \text{Norm}(0, 1)$.

s.8.1.20. Take $g(x, y) = (x, w) = (x, (x + y)/2)$. Then, $y = 2w - x$.

$$\frac{\partial(x, w)}{\partial(x, y)} = \begin{vmatrix} 1 & 0 \\ 1/2 & 1/2 \end{vmatrix} = 1/2, \quad (13.0.95)$$

$$f_{X, w}(x, w) = f_{X, Y}(x, y) \frac{\partial(x, y)}{\partial(x, w)} = \frac{1}{\pi(1+x^2)} \frac{1}{\pi(1+(2w-x)^2)} 2, \quad (13.0.96)$$

$$f_W(w) = \int_{-\infty}^{\infty} f_{X, w}(x, w) dx = \frac{2}{\pi^2} \int_{-\infty}^{\infty} \frac{1}{1+x^2} \frac{1}{1+(2w-x)^2} dx. \quad (13.0.97)$$

The expectation of a Cauchy distributed r.v. X is not well-defined because $E[|X|] = \infty$. As a consequence, taking the average of some outcomes (i.e. a sample) will also not give a sensible answer.

s.8.1.21.

s.8.1.22.

s.8.1.23.

s.8.1.24.

s.8.1.25.

s.8.1.26.

s.8.1.27.

s.8.1.28.

s.8.1.29. Incorrect: The support of T is $(0, 2)$ whereas the support of any beta distribution is $(0, 1)$. Hence, T does not have a beta distribution for some a, b .

Also see page 378 of the book for the distribution of the sum of two uniform distributions. This might help your intuition for this solution.

s.8.1.30. We use that the PDF integrates to 1:

$$1 = \int_0^1 \frac{1}{\beta(1, b)} (1-x)^{b-1} dx = \frac{1}{\beta(1, b)} \left[-\frac{1}{b} (1-x)^b \right]_0^1 = \frac{1}{\beta(1, b)b}.$$

Hence, $\beta(1, b) = \frac{1}{b}$.

s.8.1.31. The scaling factor $\beta(a, b)$ is a positive constant, so we may as well leave it out and maximize $x^{a-1}(1-x)^{b-1}$. Note that its derivative (to x) is given by

$$\begin{aligned} \frac{d}{dx} x^{a-1}(1-x)^{b-1} &= ((a-1)(1-x) - (b-1)x) x^{a-2}(1-x)^{b-2} \\ &= ((a-1) - (a+b-2)x) x^{a-2}(1-x)^{b-2}. \end{aligned}$$

Setting this to zero yields $x = \frac{a-1}{a+b-2}$ as the only candidate for an interior optimum. Since $a, b > 1$, we have $0 < x < 1$. If $a, b > 1$, then the PDF converges to 0 as $x \rightarrow 0$ or $x \rightarrow 1$, so then we conclude that $x = \frac{a-1}{a+b-2}$ indeed yields a maximum. (Think about this last sentence; most students do not use the information that $a, b > 1$ correctly.)

s.8.1.32. A prior is a distribution reflecting one's information or belief about a parameter before updating it with information.

It is harder than you might think, hardly any student gives a completely satisfactory answer here. Compare your solution to the definition above. If they are different, try to understand how exactly your solution was different and determine which definition is better.

A conjugate prior is a prior distribution such that the posterior distribution is in the same family of distributions.

s.8.1.33. Dirichlet distribution. The Beta distribution is a special case of the Dirichlet distribution, because binomial is a special case of multinomial. Of course, this can also be shown directly using the formula.

s.8.1.34. The prior is $p \sim \text{Beta}(1, 1)$. The posterior is $p|X = k \sim \text{Beta}(1 + k, 1 + n - k)$.

s.8.1.35. Let X denote the number of heads.

1. Your posterior is $p|X = 900 \sim \text{Beta}(910, 110)$.
2. Your friend's posterior is $p|X = 900 \sim \text{Beta}(901, 101)$.
3. The mean of your posterior is $\frac{910}{910+110} = \frac{91}{102} \approx 0.892$; the mean of your friend's posterior is $\frac{901}{901+101} = \frac{901}{1002} \approx 0.899$. The difference is small, so the effect of the prior distribution is small if you have a lot of data. This effect is known as *washing out the prior*.

s.8.1.36.

s.8.1.37.

s.8.1.38. This states that the PMF of the Beta-Binomial distribution,

$$P(X = k) = \binom{n}{k} \frac{\beta(a + k, b + n - k)}{\beta(a, b)},$$

sums to 1. To see this, we have to rewrite the beta functions in terms of binomial coefficients:

$$\begin{aligned} \frac{1}{\beta(a, b)} &= \frac{\Gamma(a + b)}{\Gamma(a)\Gamma(b)} = \frac{(a + b - 1)!}{(a - 1)!(b - 1)!} = (a + b - 1) \binom{a + b - 2}{a - 1}, \\ \frac{1}{\beta(a + k, b + n - k)} &= (a + b + n - 1) \binom{a + b + n - 2}{a + k - 1}. \end{aligned}$$

Plugging this in gives the result.

s.8.1.39. $V[X] = n/\lambda^2$, $E[X] = n/\lambda$, $\text{SCV} = 1/n$.

s.8.1.40.

s.8.1.41. We fill in $\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$ to find

$$f_Y(y) = \varphi(\sqrt{y}) y^{-1/2} = \frac{1}{\sqrt{2\pi}} y^{-1/2} e^{-(\sqrt{y})^2/2} = \frac{1}{\sqrt{2\pi}} y^{-1/2} e^{-y/2},$$

so $a = \frac{1}{2}$ and $\lambda = \frac{1}{2}$.

s.8.1.42. Incorrect: the scale parameters λ need to be the same *and* both random variables need to be independent.

s.8.1.43. The base case is $n = 1$. We have $\Gamma(1) = \int_0^\infty e^{-x} dx = 1 = 0!$, so the statement holds for $n = 1$. Now let $k \in \mathbb{N}$ be arbitrary and assume that the statement holds for $n = k$, i.e. that $\Gamma(k) = (k-1)!$. Then

$$\Gamma(k+1) = k\Gamma(k) = k(k-1)! = k! = ((k+1)-1)!, \quad (13.0.98)$$

so the statement also holds for $n = k+1$. By mathematical induction, we conclude that $\Gamma(n) = (n-1)!$ for all positive integers n .

s.8.1.44. Incorrect: It is the other way around, the Gamma distribution is the conjugate prior of the Poisson distribution. This statement doesn't make much sense, for example one would need to say for which parameter of the Gamma distribution it is the prior. In addition, the parameters of the Gamma distribution can be any positive real number, so the conjugate prior of (either parameter) of the Gamma distribution is a continuous distribution, so in particular not the Poisson distribution.

s.8.1.45. $X + Y \sim \text{Gamm}(11, 2)$ and $\frac{X}{X+Y} \sim \text{Beta}(4, 7)$.

s.8.1.46. 1. Minimum

2. Maximum

3. Median

s.8.2.1. From the hint, we first focus on a set $\{V \leq 0\} = \{1/T \leq 0\}$. Now, $1/T \leq 0 \iff T \leq 0$. And therefore $P\{V \leq 0\} = P\{T \leq 0\} = F_T(0)$.

If $v < 0$, then $1/T \leq v \leq 0 \iff 1/v \leq T \leq 0$. Therefore $F_V(v) = F_T(0) - F_T(1/v)$.

If $v > 0$, then $1/T \leq v$ when $T < 0$ or $T \geq 1/v$. Hence, $F_V(v) = F_T(0) + 1 - F_T(1/v)$.

s.8.2.2. a. I remember this: $f_{X,Y}(x,y) dx dy = f_{R,\Theta}(r,\theta) dr d\theta$. From this,

$$f_{R,\Theta}(r,\theta) = f_{X,Y}(x,y) \left| \frac{\partial(x,y)}{\partial(r,\theta)} \right|.$$

Now, since $x = r \cos \theta$ and $y = r \sin \theta$,

$$\frac{\partial(x,y)}{\partial(r,\theta)} = \begin{pmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{pmatrix},$$

which has determinant equal to r . It is given that $f_{X,Y}(x,y) = g(x^2 + y^2) = g(r^2)$. Hence,

$$f_{R,\Theta}(r,\theta) = f_{X,Y}(x,y)r = g(r^2)r,$$

with $r \geq 0, \theta \in [0, 2\pi]$. The RHS does not depend on θ . Hence, $f_{\Theta}(\theta)$ must be a constant.

b. Use the hint. Since g is a constant, $f_{R,\Theta}(r,\theta) \propto r$. Thus,

$$\int_0^1 \int_0^{2\pi} r \, dr \, d\theta = 2\pi(1/2)r^2|_0^1 = \pi.$$

So, $1/\pi$ is the normalization constant.

c. $f_{X,Y}(x,y) = \exp -x^2/\sqrt{2\pi} \exp -y^2/\sqrt{2\pi} = \exp -(x^2 + y^2)/2\pi = \exp -r^2/2\pi$. Indeed, $f_{X,Y}(x,y)$ has the form $g(x^2 + y^2)$. The rest is as in part b.

s.8.2.3. I always start with this line: $f_{T,U}(t,u)dt du = f_{X,Y}(x,y)dx dy$. Then,

$$\frac{\partial(t,u)}{\partial(x,y)} = \begin{pmatrix} 1/y & -x/y^2 \\ 1 & 0 \end{pmatrix} = x/y^2.$$

We don't need to take absolute signs in the last expression because X, Y are positive rvs. Next, $x = u, y = u/t$. With this,

$$f_{T,U}(t,u) = f_{X,Y}(x,y) \left(\frac{\partial(t,u)}{\partial(x,y)} \right)^{-1} = f_{X,Y}(u, u/t) y^2/x = f_X(u) f_Y(u/t) u/t^2.$$

b. Use part a.

$$f_T = \frac{1}{t^2} \int_0^\infty x f_X(x) f_Y(x/t) dx.$$

Since f_X and f_Y are not given explicitly, we cannot make further progress.

All and all, division of rvs is not so simple.

s.8.2.4. a.

$$f_{X,T}(x,t) = f_{X,Y}(x,y) \left| \frac{\partial(x,y)}{\partial(x,t)} \right|,$$

$$\frac{\partial(x,t)}{\partial(x,y)} = \begin{pmatrix} 1 & 0 \\ y & x \end{pmatrix} = x.$$

\Rightarrow

$$f_{X,T}(x,t) = f_{X,Y}(x,y)/x = f_{X,Y}(x, t/x)/x,$$

since $y = t/x$. Finally, for f_T , marginalize x out by integration.

b. Just do the algebra. With part a. you have the answer, so you can check.

s.8.2.6. a. I did not attempt any smart tricks. Take as transform $u = s/t$ and $v = s + t$, where I associate s to T_1 and t to T_2 . Then,

$$\frac{\partial(u, v)}{\partial(s, t)} = \begin{pmatrix} 1/t & -s/t^2 \\ 1 & 1 \end{pmatrix} = \frac{1}{t} + \frac{s}{t^2} = \frac{t+s}{t^2} = \frac{v}{t^2}$$

With a bit of algebra: $s = uv/(u+1)$ and $t = v/(u+1)$. Therefore the Jacobian becomes equal to $(u+1)^2/v$. Next,

$$\begin{aligned} f_{U,V}(u, v) &= f_{T_1, T_2}(s, t) \frac{\partial(s, t)}{\partial(u, v)} = f_{T_1}(uv/(u+1)) f_{T_2}(v/(u+1)) \frac{(u+1)^2}{v} \\ &= \lambda^2 \exp(-\lambda uv/(u+1) - \lambda v/(u+1)) \frac{(u+1)^2}{v} \\ &= \lambda^2 \exp(-\lambda v) \frac{(u+1)^2}{v}. \end{aligned}$$

This factors into one term with only v and another with only u . Hence, U and V are independent.

b. With the hint,

$$\begin{aligned} P\{T_1 < T_2\} &= \int_0^\infty P\{T_1 < T_2 | T_1 = s\} f_{T_1}(s) ds \\ &= \int_0^\infty P\{s < T_2 | T_1 = s\} \lambda_1 e^{-\lambda_1 s} ds \\ &= \lambda_1 \int_0^\infty e^{-\lambda_2 s} e^{-\lambda_1 s} ds = \frac{\lambda_1}{\lambda_1 + \lambda_2}. \end{aligned}$$

c. See the hint. Alice first has to wait for the first server to become free. The expected time in queue is $1/(\lambda_1 + \lambda_2)$. If server 1 is the first, then Alice spends a time $1/\lambda_1$ in service. Thus, the total time is

$$\frac{1}{\lambda_1 + \lambda_2} + \frac{\lambda_1}{\lambda_1 + \lambda_2} \frac{1}{\lambda_1} + \frac{\lambda_2}{\lambda_1 + \lambda_2} \frac{1}{\lambda_2} = \frac{3}{\lambda_1 + \lambda_2}.$$

s.8.2.7. Let us solve the question from first principles. At the end, I'll give the short solution based on Beta-Binomial conjugacy.

Let $f(p)$ be our prior density (In the exercise it is taken to be uniform). Then

$$P\{p > r\} = \int I_{p>r} f(p) dp = \int_r^1 f(p) dp$$

is our belief that $p > r$. For this exercise, we are interested in the relation $P\{p > r\} \geq c$. For instance, suppose we take $c = 0.95$, then we like to know which value for r achieves that $P\{p > r\} \geq c$?

We can start with one trial, i.e., $n = 1$. Then we analyze the case for $n = 2$, and so on, and hope to see a pattern. Here are the standard steps of Bayesian reasoning.

1. I want to know the density $f_1(p|N=1)$, i.e., the density of p after having seen one successful test. (Note here that I am careful about notation. We do $n=1$ trials, and then the number of successes is given by the random variable N .)
2. Now I use Bayes' rule:

$$f_1(p|N=1) = \frac{f_1(p, N=1)}{P\{N=1\}} = \frac{f_1(N=1|p)}{P\{N=1\}} f(p).$$

Here $f(p)$ acts as the prior density on p .

3. It is clear that $f_1(N=1|p) = p$, because we know that an item passes a test with probability p , when p is given.
4. Perhaps I don't need $P\{N=1\}$ if I can guess it (though see below), but here it is just for completeness' sake.

$$P\{N=1\} = \int_0^1 f(N=1|p)f(p)dp = \int_0^1 p dp = 1/2,$$

because the prior $f(p) = I_{p \in [0,1]}$, i.e., uniform on $[0,1]$, i.e., it is Beta(1,1).

5. With this, $f_1(p|N=1) = \frac{p}{1/2} I_{p \in [0,1]} = 2p I_{p \in [0,1]}$.

6. Thus, $P\{p > r | N=1\} = \int_0^1 I_{p>r} f_1(p|N=1) dp = \int_r^1 2p dp = 1 - r^2$.

Sometimes we are lucky and we don't have to compute the denominator in Bayes' formula. We did this earlier, but let's show again how this works.

$$f_1(p|N=1) = \frac{f_1(p, N=1)}{P\{N=1\}} \propto f_1(N=1|p)f(p) = p I_{0 \leq p \leq 1}.$$

Now $f_1(p|N=1)$ is a PDF, hence must integrate to 1. Thus, $\int_0^1 p dp = 1/2$, must be the normalization constant by which we have to divide to turn f_1 into a real PDF. In this case we don't save any work, but sometimes this really helps, in particular when dealing with integrals with Beta distributed rvs.

Now generalize to larger n , compute $f_2(p|N=2)$, then for $n=3$, and so on, until you see the pattern.

We can also directly use the ideas of the book. Starting with a prior Beta(1,1), after n 'wins', the distribution becomes Beta(1+n,1). Then,

$$P\{p > r\} = \frac{\Gamma(n+2)}{\Gamma(n+1)\Gamma(1)} \int_r^1 p^n dp = 1 - r^{n+1}. = (n+1)p^{n+1}|_r^1 = 1 - r^{n+1}.$$

s.8.2.8. a. By the hint and the fact that U_j is uniform on $[0,1]$, so that $1-U_j$ is also uniform, the last equality of the hint implies that $P\{1-U_j \leq 1-e^{-c}\} = P\{U_j \leq 1-e^{-c}\} = 1-e^{-c}$. But then, $X_j \sim \text{Exp}(1)$.

b. The sum of n iid exponentials is Gamma(n, λ). And so, if $S_n = \sum_{i=1}^n X_i$, then $P\{S_n \leq x\} = \int_0^x f(y) dy$, with $f(y)$ the gamma density with n and $\lambda = 1$.

Just to test my skills, I used MGFs, because I know that the MGF of a sum of iid rvs is the product of the MGF of one them. Since $e^{\log u} = u$,

$$\mathbb{E} \left[e^{-s \log U} \right] = \int_0^1 e^{-s \log u} du = \int_0^1 u^{-s} du.$$

If $s \geq 1$ this does not converge (convince yourself that you understand this). With $s < 1$,

$$\mathbb{E} \left[e^{-s \log U} \right] = \frac{1}{-s+1} u^{-s+1} \Big|_0^1 = \frac{1}{1-s}.$$

Therefore,

$$\begin{aligned} \mathbb{E} \left[e^{-s S_n} \right] &= \mathbb{E} \left[e^{-s \log U_1 - s \log U_2 - \dots - s \log U_n} \right] = \left(\mathbb{E} \left[e^{-s \log U} \right] \right)^n \\ &= \left(\frac{1}{1-s} \right)^n, \end{aligned}$$

and this is the MGF of a Gamma($n, \lambda = 1$) rvs.

s.8.2.9.

$$M_Y(s) = \mathbb{E}[\exp sY] = p \sum_{i=0}^{\infty} e^{s p i} q^i = p/(1 - q e^{s p}).$$

Now, use that $e^{s p} \approx 1 + s p$ for $p \ll 1$. (This is easier than using l' Hopital's rule as BH do in their solution). Hence, the denominator becomes $\approx 1 - (1-p)(1+s p) = p(1-s) - s p^2 \approx p(1-s)$ when $p \ll 1$ Hence,

$$M_Y(s) \approx p/(p(1-s)) = 1/(1-s).$$

In the limit $p \rightarrow 0$ the LHS converges to the RHS, which is the MGF of an exponential rv. For the rest, follow the solution of BH.

Here is another line of attack. Let us first use probability theory to find out what is $\sum_{i=0}^{\infty} q^i$ for some $|q| < 1$. Take $X \sim \text{Geo}(p)$, so that X corresponds to the number of failures (tails say) until we see a success (heads say). So, X corresponds to the number of tails until we see a heads. Now if we keep on throwing, then we know that eventually a heads will appear. Therefore $p + p q + p q^2 + \dots = 1$, that, is $p \sum_{i=0}^{\infty} q^i = 1$. But this implies that $\sum_{i=0}^{\infty} q^i = 1/p = 1/(1-q)$.

By similar reasoning, if we keep on throwing the coin until we see r heads then we know that $p^r \sum_{i=0}^{\infty} \binom{r+i-1}{r} q^i = 1$. Therefore,

$$\sum_{i=0}^{\infty} \binom{r+i-1}{r} q^i = \frac{1}{p^r} = \frac{1}{(1-q)^r}.$$

With this insight, for $X \sim \text{NBin}(p, n)$

$$\begin{aligned} M_X(s) &= p^r \sum_{i=0}^{\infty} \binom{r+i-1}{r} q^i e^{s i} = p^r \sum_{i=0}^{\infty} \binom{r+i-1}{r} (e^s q)^i \\ &= \frac{p^r}{(1-q e^s)^r} \approx \left(\frac{p}{p(1-s)} \right)^r, \end{aligned}$$

where we use again Taylor's expansion for $p \ll 1$.

s.9.1.2.

$$\begin{aligned} E[X] &= \int_0^1 (v-b) I_{b \geq \alpha v} dv = \int_0^{b/\alpha} v dv - b \int_0^{b/\alpha} dv \\ &= b^2/2\alpha^2 - b^2/\alpha = \frac{b^2}{\alpha^2} \frac{1-2\alpha}{2}. \end{aligned}$$

Clearly, $\alpha > 1/2$ to ensure that $E[X] > 0$.

s.9.1.3.

s.9.1.4. $E[X] = 1 + q E[X]$, because we have to throw at least once, and with probability q , we start again. Hence, $E[X] = 1/(1-q) = 1/p$.

s.9.1.5. Suppose the first throw is a success, then we need $r-1$ more successes, if the first throw is a failure, we are back at 'hole one'. Thus, $N_r = pN_{r-1} + q(1+N_r)$. Simplifying (and using that $p/(1-q) = 1$) gives $N_r = N_{r-1} + q/p$, which implies $N_r = rq/p$.

s.9.1.6.

$$N_r = N_{r-1} + p \cdot 1 + q(1+N_r) \implies N_r = N_{r-1}/p + 1/p \implies N_r = \sum_{i=1}^r 1/p^i. \quad (13.0.99)$$

s.9.1.7. Let X be the outcome of the die throw (note that X is a random variable) and let A be the event that the outcome is even. Then

$$E[X|A] = 2P\{X=2|A\} + 4P\{X=4|A\} + 6P\{X=6|A\} = \frac{1}{3} \cdot (2+4+6) = 4.$$

We conclude that $E[X|A] = 4$.

s.9.1.8. For 4: take $i = b-1$. Then solve for α in $\alpha(b-1) = \alpha(b-2)/2 + 1/2$, because $p_b = 1$. This gives $\alpha = 1/b$.

s.9.1.10.

s.9.1.12. We have that $E[Y - E[Y|X]] = 0$. Hence, $E[Y - E[Y|X]] E[X] = 0$. Then define $h(X) = E[Y|X]$ and apply BH.9.3.9 to see that $E[(Y - E[Y|X])h(X)] = 0$. From the definition of the covariance, $\text{Cov}[W, Z] = E[WZ] - E[W]E[Z]$, we have shown that both terms are zero.

s.9.1.13. 1. Since Adam keeps $b/2$ and does the gamble with $\alpha = b/2$, we have

$$E[X] = b/2 + \frac{1}{5} \cdot 4(b/2) + \frac{4}{5} \cdot 0 = 0.9b.$$

2. The computation is the same as in part 1., but with X instead of b :

$$E[Y|X] = X/2 + \frac{1}{5} \cdot 4(X/2) + \frac{4}{5} \cdot 0 = 0.9X.$$

Note that the result is a random variable.

3. Using Adam's law (and linearity of expectation), we conclude that:

$$E[Y] = E[E[Y|X]] = E[0.9X] = 0.9E[X] = 0.81b.$$

In general, if Adam would do this n times, the expected amount of money he has after n such gambles would be $0.9^n b$. This would be very difficult to show without Adam's law!

s.9.1.14. We have $E[X|N] = Np$, so using Adam's law (and linearity of expectation), we conclude that $E[X] = E[E[X|N]] = E[Np] = E[N]p = \lambda p$.

This is in accordance with $X \sim \text{Pois}(\lambda p)$, which was shown in the chicken-egg story.

Some students reported answers like $\lambda^2 p$. This is wrong, and can be immediately seen by checking units: the unit of λ being 1 per time.

Others wrote $E[X|N = n]np$, hence $E[X] = E[E[X|N]] = E[np] = np$.

Apparently, such students are not aware of the idea that $E[X|N]$ is a random variable. When this happens during the exam, you will score 0 points for that particular part of a question.

s.9.1.15. Incorrect: $E[X|A]$ is a number since A is an event, whereas $E[X|I_A]$ is a random variable since I_A is a random variable. A correct statement is $E[X|A] = E[X|I_A = 1]$.

s.9.1.16. Correct, if X and Y are independent, then $E[Y|X] = E[Y]$ which is a constant (formally, a degenerate random variable). Since the variance of a constant is 0, we conclude that $V[E[Y|X]] = 0$.

s.9.1.17. 1. We compute $E[X|X \geq a]$ as follows:

$$\begin{aligned} E[X|X \geq a] &= \int_0^\infty y f(y|A) dy \\ &= \int_0^\infty y \frac{\lambda e^{-\lambda y} I_{y \geq a}}{e^{-\lambda a}} dy \\ &= \lambda \int_a^\infty y e^{-\lambda(y-a)} dy \\ &= -y e^{-\lambda(y-a)} \Big|_a^\infty + \int_a^\infty e^{-\lambda(y-a)} dy \\ &= a - \frac{1}{\lambda} e^{-\lambda(y-a)} \Big|_a^\infty = a + \frac{1}{\lambda}. \end{aligned}$$

2. The result also follows from the memoryless property, which states that conditional on the event that $X \geq a$, we have that $X - a|X \geq a \sim \text{Exp}(\lambda)$.

s.9.1.18. 1. Note that $X|A \sim \text{Bin}(10, 0.5)$, so $E[X|A] = 10 \cdot 0.5 = 5$.

2. Note that $X|A^c \sim \text{Bin}(10, 0.8)$, so $E[X|A^c] = 10 \cdot 0.8 = 8$.

3. By LOTE we have $E[X] = P\{A\} E[X|A] + P\{A^c\} E[X|A^c] = 0.9 \cdot 5 + 0.1 \cdot 8 = 5.3$.

4. Note that $P\{B|A\} = 0.5^4$ and $P\{B|A^c\} = 0.8^4$. By LOTP we have

$$P\{B\} = P\{A\} P\{B|A\} + P\{A^c\} P\{B|A^c\} = 0.9 \cdot 0.5^4 + 0.1 \cdot 0.8^4 = 0.09721.$$

5. By Bayes' rule $P\{A|B\} = \frac{P\{B|A\}P\{A\}}{P\{B\}} \approx 0.57864$.

6. Note that $E[X|A, B] = 4 + 6 \cdot 0.5 = 7$ and $E[X|A^c, B] = 4 + 6 \cdot 0.8 = 8.8$. By LOTP with extra conditioning we have

$$P\{X|B\} = P\{A|B\} E[X|A, B] + P\{A^c|B\} E[X|A^c, B] \approx 7.75844.$$

7. By LOTE we have $P\{B\} E[X|B] + P\{B^c\} E[X|B^c] = E[X] = 5.3$. We know $P\{B\}$ and $E[X|B]$, so solving this for $E[X|B^c]$ yields $E[X|B^c] \approx 5.035$.

One or more students wrote the LOTE as $E[X] = \sum_Y E[X|Y] P\{Y\}$. This is wrong, as you cannot sum over a rv. This is correct: $E[X] = \sum_y E[X|Y=y] P\{Y=y\}$, so sum over the *outcomes* of a rv.

s.9.1.19. The marginal density of X is given by $f_X(x) = 2(1-x)$.

So the conditional density is given by $f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{I_{x \leq y}}{1-x}$. Hence,

$$E[Y|X=x] = \int_0^1 y \frac{I_{x \leq y}}{1-x} dy = \frac{1}{1-x} \int_x^1 y dy = \frac{1}{1-x} \left[\frac{1}{2} y^2 \right]_x^1 = \frac{\frac{1}{2}(1-x^2)}{1-x} = \frac{1}{2}(1+x).$$

We conclude that $E[Y|X] = \frac{1}{2}(1+X)$.

The marginal density of Y is given by $f_Y(y) = 2y$.

So the conditional density is given by $f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)} = \frac{I_{x \leq y}}{y}$. So

$$E[X|Y=y] = \int_0^1 x \frac{I_{x \leq y}}{y} dx = \frac{1}{y} \int_0^y x dx = \frac{1}{2}y.$$

We conclude that $E[X|Y] = \frac{1}{2}Y$.

Some students wrote for instance $E[X|Y] = y/2$. Apparently, such students are not aware of the idea that $E[X|N]$ is a random variable. When this happens during the exam, you will score 0 points for that particular part of a question.

s.9.1.20. Note that $E[X|X \geq a] \geq a > E[X|X < a]$. By LOTE:

$$\begin{aligned} E[X] &= P\{X \geq a\} E[X|X \geq a] + P\{X < a\} E[X|X < a] \\ &< P\{X \geq a\} E[X|X \geq a] + P\{X < a\} E[X|X \geq a] \\ &= E[X|X \geq a], \end{aligned}$$

where the inequality is strict since $P\{X < a\} > 0$.

s.9.1.21. With the hint,

$$E[N|X] = E[N-X|X] + E[X|X] = E[N-X] + X = \lambda(1-p) + X.$$

As a check, $E[E[N|X]] = E[\lambda(1-p) + X] = \lambda(1-p) + \lambda p = \lambda = E[N]$.

Here is straightforward computation. You should check each and every step as they are based on pattern recognition.

$$E[N|X=k] = \sum_{n=k}^{\infty} n P\{N=n|X=k\} \quad (13.0.100)$$

$$= \frac{1}{P\{X=k\}} \sum_{n=k}^{\infty} n e^{-\lambda} \frac{\lambda^n}{n!} \binom{n}{k} p^k (1-p)^{n-k} \quad (13.0.101)$$

$$= \frac{1}{P\{X=k\}} \sum_{n=k}^{\infty} n e^{-\lambda} \frac{1}{n!} \frac{n!}{k!(n-k)!} (\lambda p)^k (\lambda(1-p))^{n-k} \quad (13.0.102)$$

$$= \frac{e^{-\lambda p} (\lambda p)^k / k!}{P\{X=k\}} \sum_{n=k}^{\infty} n e^{-\lambda(1-p)} \frac{1}{(n-k)!} (\lambda(1-p))^{n-k} \quad (13.0.103)$$

$$= \sum_{n=k}^{\infty} n e^{-\lambda(1-p)} \frac{1}{(n-k)!} (\lambda(1-p))^{n-k} \quad (13.0.104)$$

$$= \sum_{n=0}^{\infty} (n+k) e^{-\lambda(1-p)} \frac{1}{n!} (\lambda(1-p))^n \quad (13.0.105)$$

$$= k + \sum_{n=0}^{\infty} n e^{-\lambda(1-p)} \frac{1}{n!} (\lambda(1-p))^n \quad (13.0.106)$$

$$= k + \lambda(1-p). \quad (13.0.107)$$

Hence, $E[N|X] = \lambda(1-p) + X$. Since $E[X] = \lambda p$, we get $E[N] = \lambda$ with Adam's law, as above.

s.9.1.22.

s.9.1.23. By Eve's law,

$$V[Y] = E[V[Y|X]] + V[E[Y|X]] \geq V[E[Y|X]], \quad (13.0.108)$$

since $V[Y|X] \geq 0$ for all X , which implies that $E[V[Y|X]] \geq 0$.

s.9.1.24. Conditional on Z , Y is a constant, and the variance of a constant is 0. Hence, $V[Y|Z] = 0$.

s.9.1.25. Incorrect. Counterexample: Let $Y \sim \text{Bern}(1/2)$ and A be the event $Y = 0$. then $\text{Var}(Y|A)$ and $\text{Var}(Y|A^c)$ are both 0, but $\text{Var}(Y) = 1/4$.

s.9.1.26. $V[Y|X]$ is a random variable, but $V[Y|X=x]$ is a constant.

s.9.1.27. Define $g(X) = E[Y|X]$. Then,

$$E[(Y - E[Y|X])^2|X] = E[(Y - g(X))^2|X] \quad (13.0.109)$$

$$= E[Y^2 - 2Yg(X) + g(X)^2|X] \quad (13.0.110)$$

$$= E[Y^2|X] - 2E[Yg(X)|X] + E[g(X)^2|X] \quad (13.0.111)$$

$$= E[Y^2|X] - 2g(X)E[Y|X] + g(X)^2 \quad (13.0.112)$$

$$= E[Y^2|X] - 2g(X)^2 + g(X)^2 \quad (13.0.113)$$

$$= E[Y^2|X] - (E[Y|X])^2 \quad (13.0.114)$$

s.9.1.28. Using Eve's Law we have

$$V[W] = V[E[W|X]] + E[V[W|X]] = V[0] + E[X^2] = 0 + \mu^2 + \sigma^2 = \mu^2 + \sigma^2. \quad (13.0.115)$$

s.9.2.1. Using the hint, for a. $E[T] = E[\sum_j E[R_j] I_{R=j}] = \sum_j E[R_j] E[I_{R=j}] = (\mu_1 + \mu_2 + \mu_3)/3$.

For b., $E[T^2|R]$, realize that $E[R_j^2] = \mu_j^2 + \sigma_j^2$, because $V[R_j] = E[R_j^2] - (E[R_j])^2$. Finally, with these ideas,

$$\begin{aligned} V[E[T|R]] &= E[(E[T|R])^2] - (E[E[T|R]])^2 \\ &= (\mu_1^2 + \mu_2^2 + \mu_3^2)/3 + (\mu_1 + \mu_2 + \mu_3)^2/9. \end{aligned}$$

s.9.2.2.

$$\begin{aligned} E[X_{n+1}|X_n = 100] &= 100 + pf100 - (1-p)f100 = 100(1-f+2pf) \\ E[X_{n+1}|X_n] &= X_n(1-f+2pf) \\ E[X_{n+1}] &= (1-f+2pf)E[X_n] \\ E[X_{n+1}] &= (1-f+2pf)^2 E[X_{n-1}] = (1-f+2pf)^{n+1} X_0. \end{aligned}$$

s.9.2.4. Use that $P^2 = P$ (indicator function), Adam and Eve, and that $N \sim \text{Pois}(8\lambda)$,

$$\begin{aligned} E[Y|P] &= E[PX|P] = E[X] E[P|P] = \mu P, & V[Y|P] &= V[XP|P] = P^2 V[X|P] = P\sigma^2 \\ E[Y] &= \mu p, & V[Y] &= E[V[Y|P]] + V[E[Y|P]] = \sigma^2 p + \mu^2 p(1-p), \\ E[S|N] &= N E[Y], & V[S|N] &= N V[Y] \\ E[N] &= 8\lambda, & V[N] &= 8\lambda. \end{aligned}$$

Now use BH.9.6.1. It's just a matter of filling in how.

s.9.2.5. a. Here you should assume that the X_i are not yet known. Thus, the expectation over X_i is taken with respect to the CDF F_X . Using the independence of X_j and S_j , $I_{S_j=i} I_{S_j=k} = 0$ if $i \neq k$, and that $E[I_{S_j=k}] = 1/n$,

$$\begin{aligned} E[Y_j] &= \sum_i E[X_i] E[I_{S_j=i}] = \mu, \\ E[Y_j^2] &= E\left[\sum_k \sum_l X_k X_l I_{S_j=k} I_{S_j=l}\right] = E\left[\sum_k X_k^2 I_{S_j=k}\right] = \sum_k E[X_k^2] n^{-1} = E[X^2], \\ V[Y_j] &= E[Y_j^2] - (E[Y_j])^2 = \sigma^2. \end{aligned}$$

b. Now we are given the outcomes (samples) $X_i = x_i$ of n experiments. I prefer to write $D = X_1, \dots, X_n$ as it is shorter. Noting that S_j and D are independent, and that $E[X_k|D] = X_k$,

$$E[Y_j|D] = \sum_k X_k E[I_{S_j=k}|D] = \frac{1}{n} \sum_k X_k := \bar{X}$$

Observe that this average need not be the same as μ !

The conditional variance. Since S_j and S_k are independent when $j \neq k$, it must be that $Y_j|D$ and $Y_k|D$ are also conditionally independent. Moreover, $\{Y_j|D\}$ are conditionally iid. Therefore,

$$\begin{aligned} E[Y_j^2|D] &= E\left[\sum_k \sum_l X_k X_l I_{S_j=k} I_{S_j=l} | D\right] \\ &= E\left[\sum_k X_k^2 I_{S_j=k} | D\right] = \sum_k X_k^2 E[I_{S_j=k} | D] \\ &= \frac{1}{n} \sum_k X_k^2, \\ V[Y_j|D] &= \frac{1}{n} \sum_k X_k^2 - (\bar{X})^2 = \frac{1}{n} \sum_k (X_k - \bar{X})^2 = \frac{n-1}{n} \sigma^2, \\ V[\bar{Y}|D] &= V\left[\frac{1}{n} \sum_j Y_j | D\right] = \frac{1}{n^2} \sum_j V[Y_j | D] = \frac{1}{n} V[Y_1 | D]. \end{aligned}$$

c. For $E[\bar{Y}]$ use linearity and Adam's law:

$$E[\bar{Y}] = E[E[\bar{Y}|D]] = \frac{1}{n} \sum_k E[X_k] = E[X] = \mu.$$

Here are the details for $V[\bar{Y}]$. Using BH.6.3.3 and BH.6.3.4,

$$\begin{aligned} E[V[\bar{Y}|D]] &= \frac{1}{n} E[V[Y_1|D]] = \frac{1}{n^2} E\left[\sum_{i=1}^n (X_i - \bar{X})^2\right] \\ &= \frac{n-1}{n^2} E\left[\frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2\right] = \frac{n-1}{n^2} E[S_n^2] = \frac{(n-1)\sigma^2}{n^2} \\ V[E[\bar{Y}|D]] &= V[\bar{X}] = \frac{1}{n^2} \sum_i V[X_i] = \frac{1}{n} \sigma^2. \end{aligned}$$

Now use Eve's law to add both terms to get $V\bar{Y}$.

d. We add randomness twice, first we draw samples to get D , and then we draw randomly from D .

The extra exercise: immediate from Example 1.4.22. We are not interested in the sequence of the bootstrap sample. BTW, the story that goes for me with this example is the 'balls and bars story'. I have n balls to distribute over k boxes. Hence, there are $k-1$ bars to separate the boxes. For the bootstrap sample, I have to distribute n bootstrap samples (the X_i^*) over n boxes (the initial sample X_i .)

If n is small, say $n = 4$. Does it make sense to take more than 1000 bootstrap samples?

s.9.2.7. a. Using the hint gives us $E[N|\lambda] = \lambda$ and $V[N|\lambda] = \lambda$.

Now use Adam and Eve.

b. Just copy the formulas of BH.9.6.1

- c. With the hint, observe that $\text{Exp}((1)) = \Gamma(1, 1)$. In the relevant formula of BH.8.4.5 ($P\{Y = y\}$), take $t = r_0 = b_0 = 1$ and conclude that $P\{N = n\} = 2^{-n-1}$. Hence, $N \sim \text{Geo}(1/2)$.
- d. Same story. The relevant formula is $f_1(\lambda|y)$.

s.9.2.10. a. From the hint,

$$\begin{aligned} E[T] &= E[E[T|p]] = \frac{1}{\beta(a, b)} \int_0^1 \frac{1}{p} p^{a-1} (1-p)^{b-1} dp \\ &= \frac{1}{\beta(a, b)} \int_0^1 p^{a-2} (1-p)^{b-1} dp = \frac{\beta(a-1, b)}{\beta(a, b)} \\ &= \frac{a+b-1}{a-1} = 1 + \frac{b}{a-1}. \end{aligned}$$

To get the last equation, use the definition of $\beta(a, b)$ in terms of factorials (see the Bayes' billiards story) to simplify. This is easy, many terms cancel.

b. Take $Y = 1 + G$, then Y has the first success distribution since G is geometric. Hence, $E[Y] = (a+b)/a = 1 + b/a$. Clearly, this is smaller than $1 + b/(a-1) = E[T]$.

But why is this so?

I must miss something here. The prior is $\text{Beta}(a, b)$. Then Beta-Binomial conjugacy story, we assume that Vishy won $a-1$ games, and lost $b-1$ games. My guess for Vishy winning the next game would be $(a-1)/(a+b-2)$, not $a/(a+b)$. But I make an error here. Check the BH problem 9.57. You'll see that we should indeed use $a/(a+b)$! Tricky!

c. Immediate from BH.8.3.3: $p|X = 7 \sim \text{Beta}(a+7, b+3)$.

s.9.2.11. a. By the hint,

$$f(p|X_1 = x_1) \propto f(p, X_1 = x_1) \propto p^{a-1} q^{b-1} p^{x_1} q^{1-x_1} \propto p^{a+x_1-1} q^{b+(1-x_1)-1}.$$

Hence, $p|X_1 = x_1 \sim \text{Beta}(a+x_1, b+(1-x_1))$. We can now use this as prior to see that $p|X_1 = x_1, X_2 = x_2 \sim \text{Beta}(1+x_1+x_2, 1+(1-x_1)+(1-x_2))$, and so on. Hence, $p|X_1, \dots, X_n \sim \text{Beta}(1+S_k, 1+n-S_k)$.

b. With the hint, $P\{X_{n+1} = 1|p\} = p$ and $P\{S_n = k|p\} = \binom{n}{k} p^k q^{n-k} \propto p^k q^{n-k}$. Also $X_{n+1}|p$ and $S_n|p$ are conditionally independent. Therefore,

$$P\{X_{n+1} = 1, S_n = k|p\} \propto p p^k q^{n-k} = p^{k+1} q^{n-k},$$

which in turn implies that

$$P\{X_{n+1} = 1|S_n = k, p\} \propto p^{k+1} q^{n-k}.$$

Hence, $X_{n+1}|S_n = k, p \sim \text{Beta}(k+2, n-k+1)$. Now, $X_{n+1} \in \{0, 1\}$, so that $P\{X_{n+1} = 1|S_n = k, p\} = E[X_{n+1}|S_n = k, p] = (k+2)/(n+3)$, since $X_{n+1}|S_n = k, p \sim \text{Beta}(k+2, n-k+1)$.

The last step is to realize that $E[X_{n+1}|S_n = k] = E[E[X_{n+1}|S_n = k, p]|S_n = k]$.

Here is another way to get the same result.

$$\begin{aligned}
 P\{S_n = k\} &= \frac{1}{n+1}, \text{ by Bayes' billiard,} \\
 P\{X_{n+1} = 1, S_n = k\} &= \int_0^1 P\{X_{n+1} = 1, S_n = k | p\} f(p) dp = \int_0^1 p \binom{n}{k} p^k (1-p)^{n-k} f(p) dp \\
 &= \frac{k+1}{n+1} \int_0^1 \binom{n+1}{k+1} p^{k+1} (1-p)^{n-k} f(p) dp \\
 &= \frac{k+1}{n+1} \frac{1}{n+2}, \text{ again with Bayes' billiard,} \\
 P\{X_{n+1} = 1 | S_n = k\} &= P\{X_{n+1} = 1, S_n = k\} / P\{S_n = k\}.
 \end{aligned}$$

Now simplify.

s.9.2.12.

$$\begin{aligned}
 E[S_n | p] &= np \\
 E[p] &= \frac{1}{\beta(a, b)} \int_0^1 p p^{a-1} q^{b-1} dp = \frac{\beta(a+1, b)}{\beta(a, b)} = \frac{a}{a+b} = 1/2 \\
 E[E[S_n | p]] &= n E[p] = n/2. \\
 V[S_n | p] &= npq \\
 E[V[S_n | p]] &= n E[pq] = n E[p] - n E[p^2] = n/2 - n E[p^2] \\
 E[p^2] &= \frac{1}{\beta(a, b)} \int_0^1 p^2 p^{a-1} q^{b-1} dp = \frac{\beta(a+2, b)}{\beta(a, b)} = \frac{a(a+1)}{(a+b)(a+b+1)} = \frac{2}{2 \cdot 3} = 1/3 \\
 V[E[S_n | p]] &= V[np] = n^2 V[p] = n^2/12.
 \end{aligned}$$

The rest of Eve's law is now trivial.

b. We start with a Beta(1, 1) prior on p . After the first win, the prior gets updated to Beta(1 + 1, 1), after a loss to Beta(1, 1 + 1). Reasoning like this, after a wins and $j - a$ losses, the distribution for a win becomes Beta(1 + a , $a + j - 1$). Therefore, by using the hint in the book, $E[p | S_j = a] = (a + 1)/(j + 2)$.

c. When somebody doesn't give me any information about what team can win, then any outcome must be equally likely. (What else can it be?) This is also my way to understand the expression in BH.8.3.2. Hence, $P\{X = k\} = 1/(n + 1)$. Observe that we use the prior $p \sim \text{Beta}(1, 1)$.

When the prior is Beta($a, j - a$), we should get the negative hypergeometric distribution, see the remark in BH.8.3.3.

d. Shanille scores the first and missed the second. Hence, there are 98 shots left, out which she has to score 49. Thus, we ask for $P\{S_{98} = 49 | p\}$, where $p \sim \text{Beta}(a = 1, b = 1)$ is the prior since she hit $a = 1$ out of $a + b = 2$ shots. This places us in the situation of part c above, with $n = 98$. Hence, $P\{S_{98} = 49 | p\} = 1/99$.

s.10.1.1. f is the expectation of something non-negative. Then, work out the square and apply linearity of the expectation. As $f \geq 0$, it can have most one root, hence $D \leq 0$. But $D = 4E[XY] - 4E[X^2]E[Y^2]$.

s.10.1.3.

s.10.1.4.

s.10.1.5.

s.10.1.6.

s.10.1.8. Equality holds for functions that are both convex and concave. The only functions that are both convex and concave are affine functions, i.e., functions of the type $g(x) = ax + b$. Assuming that g is twice differentiable, we can show this as follows. Convexity is equivalent to $g''(x) \geq 0$ and concavity is equivalent to $g''(x) \leq 0$. This means $g''(x) = 0$ and the only functions for which this holds are affine functions.

If you like maths, consider generalizing the condition. Is it necessary to assume that g is twice differentiable? For instance, it is not hard to prove that a convex function is continuous. Consider now a point at which g is convex and concave at the same time, does it follow that g is twice differentiable at such a point?

s.10.1.10. In the equation of the hint, take expectations at both sides. Realize that $E[I_{X \geq a}] = P\{X \geq a\}$. Next, for any rv, $|X| \geq 0$. Hence, we can apply the simple form of Markov's inequality to get the result of the book.

s.10.1.11. By definition, equation (1) is Chebyshev's inequality. Letting $a = c\sigma_X$ we get (4). Equation (3) follows from multiplying (1) by -1 , adding 1 and using the complement rule. Equation (2) is not equivalent to any of the others, as this is not how reversing inequalities works.

s.10.1.12. This will result in the trivial bound $P\{|X - \mu| \geq a\} \leq B$, for some $B \geq 1$. But we already know that every probability is at most one. So the bound does not tell us anything interesting.

s.10.1.13. In this (pathological) example we get from Markov's inequality that $P(X \geq 2) \leq \frac{E(X)}{2} = \frac{1}{4}$. This means the Markov bound is tight, as it is equal to the probability that X exceeds 2. From Chernoff's bound we get

$$P(X \geq 2) \leq \frac{E(e^{tX})}{e^{2t}} = \frac{3 + e^{2t}}{4e^{2t}} = \frac{1}{4} \left(1 + \frac{3}{e^{2t}} \right) > \frac{1}{4} \quad \forall t > 0.$$

Hence here the Markov bound is tighter. We use the facts from probability theory that $E(X) = \frac{1}{2}$ and that $E(e^{tX}) = \frac{3}{4} + e^{2t}\frac{1}{4}$ in this example.

s.10.1.14.

s.10.1.15.

s.10.1.16.

s.10.1.17. Divide by the std corresponds to the standard transformation $(X - \mu)/\sigma$. Like this, I don't have to remember anything new. Algebra gives the formula of the book.

s.10.1.18.

s.10.1.19.

s.10.1.20.

s.10.1.21. Here is the reason. \bar{Z}_n is the sum of n normal rvs Z_j , hence normal itself. As each of these Z_j is standard normal, $E[\bar{Z}_n] = 0$, and $V[\bar{Z}_n] = n^{-2} \sum_j V[Z_j] = 1/n$, by independence. Therefore, $\sqrt{n}\bar{Z}_n \sim N(0, 1) \implies (\sqrt{n}\bar{Z}_n)^2 \sim \chi_1^2$, where we use Definition 10.4.1 and Theorem 10.4.2 in the last step.

s.10.1.22.

s.10.1.23. Let X be the r.v. corresponding to the number of heads. Then $X \sim \text{Binomial}(100, \frac{1}{2})$, which has moments $E[X] = 100 \cdot \frac{1}{2} = 50$ and $V[X] = 100 \cdot \frac{1}{2} \cdot \frac{1}{2} = 25$. By symmetry of the $\text{Binomial}(100, \frac{1}{2})$ distribution,

$$P\{X \leq 40\} = P\{X \geq 60\}. \quad (13.0.116)$$

Hence, using Chebyshev's inequality,

$$P\{X \leq 40\} = \frac{1}{2} P\{|X - 50| \geq 10\} \quad (13.0.117)$$

$$= \frac{1}{2} P\{|X - 50| \geq 10\} \quad (13.0.118)$$

$$\leq \frac{1}{2} \frac{V[X]}{10^2} \quad (13.0.119)$$

$$= \frac{1}{2} \frac{25}{100} = \frac{1}{8}. \quad (13.0.120)$$

Hence,

$$P\{X \leq 40\} \leq \frac{1}{8}. \quad (13.0.121)$$

s.10.1.25. We assemble m observations of Y_j (hence, we throw the coin nm times). Suppose we see M times that $|Y_j - \mu| > \epsilon$. Then we expect that $M/m < \sigma^2/n\epsilon$.

Thus, Chebyshev's inequality makes a statement about sample means of size n , say.

s.10.1.26. First fix some $\epsilon > 0$. Now take some n and determine the fraction of outliers, that is, count how many of the sample means $Y_1 = \sum_{i=1}^n X_i/n, Y_2 = \sum_{i=n+1}^{2n} X_i/n, \dots$ lie outside the interval $[\mu - \epsilon, \mu + \epsilon]$ and divide by the number of samples taken. The WLLN says this: If the sample averages Y_1, Y_2 are taken over larger sets of the X_j , i.e., n is larger so that we put more throws in a batch, then the fraction of outliers become smaller.

s.10.1.27. The SLLN says nothing about individual sample paths, i.e., strings of outcomes like H, T, H, T, \dots . In fact, the probability of obtaining any particular sample path has zero probability. Instead, the SLLN makes a statement about sets of sample paths. For the coin it says that it is virtually impossible to pick a path from the set of paths whose long-run fraction of heads is not equal to $1/2$.

s.10.2.3. Take W as in the hint and $Z = 1$. By the inequality of Cauchy-Schwarz, $(E[W])^2 \geq E[W^2]$. The LHS is σ^4 , the RHS is $E[(X - \mu)^4]$. The rest follows right away from the definition of kurtosis.

s.10.2.4. a. \leq Immediate from the hint.

b. $=$: immediate from the hint

c.

$$P\{X > Y - 3\} = P\{X > Y + 3\} + P\{Y - 3 \leq X \leq Y + 3\}.$$

Both terms on the RHS are non-negative.

d. Use the hint. $(E[XY])^2 \leq E[X^2] E[Y^2] = (E[X^2])^2 \leq E[X^4]$, where we use that X and Y are iid, so that $E[X^2]$ and $E[Y^2]$ are equal.

e. $=$: since X and Y are independent, $V[Y|X] = V[Y]$.

f. From the hint, $P\{|X + Y| > 3\} \leq E[|X + Y|]/3 \leq E[|X|]/3 + E[|Y|]/3 = 2E[|X|]/3 \leq E[|X|]$. Why is there not an $<$ in the last step?

s.10.2.6. a. I did things a bit differently than in the book. Take $S_n = \sum_{i=1}^n X_i$ with $X_i \sim \text{Bern}(p)$. Then I know this:

$$P\{S_n = k\} = \binom{n}{k} p^k (1-p)^{n-k} \rightarrow e^{-\lambda} \lambda^k / k! = P\{N = k\}, \quad \text{if } N \sim \text{Pois}(\lambda),$$

for $n \rightarrow \infty, p \rightarrow 0$ but such that $pn = \lambda$. I also know from the CTL that $S_n \sim N(np, np(1-p))$ if n becomes large. But, $N(np, np(1-p)) \rightarrow N(\lambda, \lambda)$ in the above limit. Now take $\lambda = n$ to see that $\text{Pois}(\lambda) \sim N(n, n)$.

b. Check the solution manual. Then, with $\mu = \sigma = \lambda = n$, and $n \gg 1$,

$$\begin{aligned} \Phi(n + 1/2) - \Phi(n - 1/2) &= \frac{1}{\sqrt{2\pi\sigma^2}} \int_{n-1/2}^{n+1/2} e^{-(x-\mu)/2\sigma^2} dx \\ &= \frac{1}{\sqrt{2\pi n}} \int_{-1/2}^{1/2} e^{-x^2/2n} dx \\ &= \frac{1}{\sqrt{2\pi n}} \int_{-1/2}^{1/2} (1 - x^2/2n) dx \\ &= \frac{1}{\sqrt{2\pi n}} (1 - 1/(24n)). \end{aligned}$$

So, we found another term to approximate $n!$ yet better.

s.10.2.7. Since $X_n \sim \text{Pois}(n)$, $E[X_n] = n$, $V[X_n] = n$. Using the hints, with Y_n the standardized version of X_n :

$$\begin{aligned} M_{Y_n}(s) &= \sum_{i=0}^{\infty} e^{-n} n^i / i! \cdot e^{s(i-n)/\sqrt{n}} = e^{-n} e^{s\sqrt{n}} \sum_{i=0}^{\infty} (n e^{s/\sqrt{n}})^i / i! \\ &= \exp\{-n + s\sqrt{n} + n e^{-s/\sqrt{n}}\}. \end{aligned}$$

With Taylor's expansion for e^x to second order,

$$-n + s\sqrt{n} + n e^{-s/\sqrt{n}} \approx -n + s\sqrt{n} + n(1 - s/\sqrt{n} + s^2/2n) = s^2/2.$$

Now follow the proof of the CTL, BH.10.3.1.

s.10.2.8. a. Define I_n as the success indicator: it is 1 if I win, and 0 if I loose. For round 1, suppose I win, then $Y_1 = Y_0/2 + 1.7Y_0/2 = 1.35Y_0$. If I loose, $Y_1 = Y_0/2 + 0.5Y_0/2 = 0.75Y_0$. Therefore,

$$Y_n = Y_{n-1}(1.35)^{I_n}(0.75)^{1-I_n}.$$

With this expression, the rest is simple, just follow BH.10.3.7. It turns out that $Y_n \rightarrow \infty$ as $n \rightarrow \infty$.

b. Use the hint.

$$\begin{aligned} Y_n &= Y_{n-1}(1 + 0.7\alpha)^{I_n}(1 - 0.5\alpha)^{1-I_n} \implies \\ \log Y_n &= \log Y_{n-1} + I_n \log(1 + 0.7\alpha) + (1 - I_n) \log(1 - 0.5\alpha) \\ &= \log Y_0 + \log(1 + 0.7\alpha) \sum_{i=1}^n I_i + \log(1 - 0.5\alpha) \cdot \sum_{i=1}^n (1 - I_i) \end{aligned}$$

By the strong law, $\sum I_i/n \rightarrow 1/2$ and $\sum(1 - I_i)/n \rightarrow 1/2$. Therefore

$$n^{-1} \log Y_n \rightarrow 0.5 \log(1 + 0.7\alpha) + 0.5 \log(1 - 0.5\alpha) = 0.5 \log((1 + 0.7\alpha)(1 - 0.5\alpha)) = g(\alpha)$$

For the maximum, take the derivative with respect to α . This gives $\alpha = 2/7$.

s.10.2.10. a. $P\{N = n\} = P\{X_1 < 1, X_2 < 1, \dots, X_{n-1} < 1, X_n > 1\}$. But, then N must have the first success distribution, and $N - 1$ be geometric.

b. Let X_i be the inter-arrival time between jobs $i-1$ and i . Then $S_n = \sum_{i=1}^n X_i$ is the arrival time of job n . We want that $S_{M-1} < 10 \leq S_M$. Since the X_i are $\sim \text{Exp}(\lambda)$, $S_n \sim \text{Pois}(\lambda t)$.

c. The sum of n iid $\text{Exp}(1)$ rvs is $\text{Gamma}(n, 1)$. Since \bar{X}_n has mean 1, $X_n \sim \text{Gamma}(n, n)$. Then $V[X_n] = 1/n$ (I just looked it up in the back of the book). By the CLT, \bar{X}_n is approximated well by a $\text{Norm}(\mu, \sigma^2)$ rv with $\mu = 1, \sigma^2 = 1/n$.

s.11.0.1. I have remembered that $V[X] = E[X]$ when $X \sim \text{Exp}(\lambda)$. Since $V[X] = E[X^2] - (E[X])^2$, $E[X^2] = 2/\lambda^2$. Applying this to L , we see that $E[L^2] = 2/(2\lambda)^2 = 1/2\lambda^2$. Moreover, $V[L] = 1/4\lambda^2$.

Next, $f_M(x) = 2f_x(x)F_Y(x)$. Hence,

$$E[M^2] = \int x^2 2\lambda e^{-\lambda x} (1 - e^{-\lambda x}) dx = 2 \int x^2 \lambda e^{-\lambda x} dx - \int x^2 2\lambda e^{-2\lambda x} dx.$$

The first integral is just 2 times $E[X^2]$, the second is $E[L^2]$. Hence, $E[M^2] = 4/\lambda^2 - 1/2\lambda^2 = 7/2\lambda^2$. Finally, $V[M] = 7/2\lambda^2 - 9/4\lambda^2 = 5/4\lambda^2$.

s.11.0.2. 1. Since (X, Y) are bivariate normally distributed, every linear combination of X and Y is normally distributed. Note that every linear combination of $(X + Y)$ and $(X - Y)$ can be written as a linear combination of X and Y . Hence, every linear combination of $(X + Y)$ and $(X - Y)$ is normally distributed. Hence, $(X + Y, X - Y)$ is bivariate normally distributed.

2. By the story above, both X and Y are normally distributed. We have

$$E[X + Y] = E[X] + E[Y] = \mu + \mu = 2\mu, \quad (13.0.122)$$

and

$$E[X - Y] = E[X] - E[Y] = \mu - \mu = 0. \quad (13.0.123)$$

Moreover,

$$V[X + Y] = V[X] + V[Y] + 2\text{Cov}[X, Y] = 2\sigma^2 + 2\rho\sigma^2 = 2(1 + \rho)\sigma^2. \quad (13.0.124)$$

Similarly,

$$V[X - Y] = V[X] + V[-Y] + 2\text{Cov}[X, -Y] = V[X] + V[Y] - 2\text{Cov}[X, Y] \quad (13.0.125)$$

$$= 2\sigma^2 - 2\rho\sigma^2 = 2(1 - \rho)\sigma^2. \quad (13.0.126)$$

So we have found that $X + Y \sim N(2\mu, 2(1 + \rho)\sigma^2)$ and $X - Y \sim N(0, 2(1 - \rho)\sigma^2)$.

3. We have

$$\text{Cov}[X + Y, X - Y] = \text{Cov}[X, X] - \text{Cov}[X, Y] + \text{Cov}[Y, X] - \text{Cov}[Y, Y] \quad (13.0.127)$$

$$= V[X] - V[Y] = \sigma^2 - \sigma^2 = 0. \quad (13.0.128)$$

Write $U = X + Y$, $V = X - Y$. Plugging all the parameters into the formula for the joint pdf of a bivariate normal distribution (see https://en.wikipedia.org/wiki/Multivariate_normal_distribution#Bivariate_case), we obtain

$$f_{U,V}(u, v) = \frac{1}{2\pi\sqrt{2(1+\rho)\sigma^2}2(1-\rho)\sigma^2} \exp\left(-\frac{1}{2}\left[\frac{(u-2\mu)^2}{2(1+\rho)\sigma^2} + \frac{v^2}{2(1-\rho)\sigma^2}\right]\right). \quad (13.0.129)$$

s.11.0.3. This question tests your modeling skills too.

In hindsight, the questions have to be reorganized a bit. The capital at the end of the i th week is $I_i = I_{i-1} + X_i - C_i$.

Suppose claims arrive at the beginning of the week, and contributions arrive at the end of the week (people prefer to send in their claims early, but they prefer to pay their contribution as late as possible). If we don't have sufficient money in cash, then we cannot pay a claim. Thus, $\max\{I_0 - C_1\}$ is our capital just before the contribution arrives. Hence, I'_1 is our capital at the end of week 1 under the assumption that we never pay out more than we have in cash. Likewise for I'_2 .

\bar{I}_n is the lowest capital we have seen for the first n weeks.

In the supermarket setting, I_i is our inventory; we can be temporarily out of stock, but as soon as new deliveries—so called replenishments—arrive then we serve the waiting customers immediately. The model with I' corresponds to a setting in which we consider unmet demand as lost.

$$P\{I_0 \leq 0\} = P\{2 + X_1 - C_1 < 0\} = \frac{1}{10} \sum_{i=1}^{10} P\{C_1 > 2 + i\} = \frac{1}{10} \sum_{i=1}^5 P\{C_1 > 2 + i\} \quad (13.0.130)$$

$$= \frac{1}{10} \sum_{i=1}^5 \frac{6-i}{9}. \quad (13.0.131)$$

When grading, I realized that question 8 was not quite reasonable to ask as an exam question. We graded this leniently. As I find it too boring to compute these probabilities by hand, here is the python code. The ideas in the code are highly interesting and useful. The main data structure here is a dictionary, one of the most used data structures in python. I don't have the R code yet, so if you take the (unwise) decision to stick to only R, you have to wait a bit until somebody sends me the R code for this problem.

Python Code

```

1  C = {}
2  for i in range(0, 9):
3      C[i] = 1 / 9
4
5  X = {}
6  for i in range(1, 11):
7      X[i] = 1 / 10
8
9
10 I0 = 2
11
12 I1 = {}
13 for k, p in X.items():
14     for l, q in C.items():
15         i = I0 + k - l
16         I1[i] = I1.get(i, 0) + p * q

```

```

17
18 print("I1, ", sum(I1.values())) # check
19
20
21 # compute P(I1<0):
22 P = sum(r for i, r in I1.items() if i < 0)
23 print(P)
24
25
26 I2 = {}
27 for i, r in I1.items():
28     for k, p in X.items():
29         for l, q in C.items():
30             j = i + k - l
31             I2[j] = I2.get(j, 0) + r * p * q
32
33 print("I2 ", sum(I2.values())) # just a check
34
35 # compute P(I2<0):
36 P = sum(r for i, r in I2.items() if i < 0)

```

Interestingly, $I'_i \geq 1$. (This is so simple to see that I first did it wrong.)

Mistake: note that X_i and C_i are discrete rvs, not continuous. The sum of two uniform random variables is not uniform. For example, think of the sum of two die throws. Is getting 2 just as likely as getting 7?

s.11.0.4. We have

$$\text{Cov}[X, Y] = \text{Cov}[X, X^2] = E[XX^2] - E[X]E[X^2] = 0 - 0 \cdot 2.5 = 0. \quad (13.0.132)$$

Hence, $\text{Corr}(X, Y) = 0$.

Yes, for instance, take $X \sim \text{Unif}(\{0, 1\})$. Then,

$$\text{Cov}[X, Y] = E[XX^2] - E[X]E[X^2] = 0.5 - 0.5 \cdot 0.5 = 0.25. \quad (13.0.133)$$

s.11.0.5. 1. The interpretation is: the time until the first component fails. That is, the time until the machine stops working.

2. Let $\lambda = 10$. We have

$$P\{\text{machine not failed at time } T\} = P\{\min\{X_1, X_2\} > T\} \quad (13.0.134)$$

$$= P\{X_1 > T, X_2 > T\} \quad (13.0.135)$$

$$= P\{X_1 > T\} P\{X_2 > T\} \quad (13.0.136)$$

$$= e^{-\lambda T} \cdot e^{-\lambda T} \quad (13.0.137)$$

$$= e^{-(2\lambda)T} \quad (13.0.138)$$

$$= e^{-20T} \quad (13.0.139)$$

$$(13.0.140)$$

3. Note that

$$P\{\min\{X_1, X_2\} \leq T\} = 1 - P\{\min\{X_1, X_2\} > T\} = 1 - e^{-20T}. \quad (13.0.141)$$

Note that this is the cdf of an exponential distribution with parameter 20. Hence, $\min\{X_1, X_2\} \sim \exp(20)$.

4. The expected time until the machine fails is

$$E[\min\{X_1, X_2\}] = 1/20, \quad (13.0.142)$$

i.e., 3 minutes. Apparently, the machine is not very robust.

s.11.0.6. 1. We have

$$P\{X + Y > 1\} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I_{X+Y>1} f_{X,Y}(x, y) dy dx \quad (13.0.143)$$

$$= \int_0^1 \int_{1-x}^1 \frac{6}{7} (x+y)^2 dy dx \quad (13.0.144)$$

$$= \frac{6}{7} \int_0^1 \left[\frac{1}{3} (x+y)^3 \right]_{y=1-x}^1 dx \quad (13.0.145)$$

$$= \frac{2}{7} \int_0^1 \left((x+1)^3 - (x+1-x)^3 \right) dx \quad (13.0.146)$$

$$= \frac{2}{7} \int_0^1 \left((x+1)^3 - 1 \right) dx \quad (13.0.147)$$

$$= \frac{2}{7} \left[\frac{1}{4} (x+1)^4 - x \right]_{x=0}^1 \quad (13.0.148)$$

$$= \frac{1}{14} \left[(x+1)^4 - 4x \right]_{x=0}^1 \quad (13.0.149)$$

$$= \frac{1}{14} \left(((1+1)^4 - 4) - ((0+1)^4 - 0) \right) \quad (13.0.150)$$

$$= \frac{1}{14} (16 - 4 - 1) \quad (13.0.151)$$

$$= \frac{11}{14}. \quad (13.0.152)$$

2. We have

$$\text{Cov}[U, V] = E[UV] - E[U]E[V]. \quad (13.0.153)$$

First, we compute

$$E[UV] = \int_0^1 \int_0^{1-u} 2uv \, dv \, du \quad (13.0.154)$$

$$= \int_0^1 [uv^2]_{v=0}^{1-u} \, du \quad (13.0.155)$$

$$= \int_0^1 (u(1-u)^2 - 0) \, du \quad (13.0.156)$$

$$= \int_0^1 u(1-2u+u^2) \, du \quad (13.0.157)$$

$$= \int_0^1 (u-2u^2+u^3) \, du \quad (13.0.158)$$

$$= \left[\frac{1}{2}u^2 - \frac{2}{3}u^3 + \frac{1}{4}u^4 \right]_{u=0}^1 \quad (13.0.159)$$

$$= \frac{1}{2} - \frac{2}{3} + \frac{1}{4} \quad (13.0.160)$$

$$= \frac{1}{12}. \quad (13.0.161)$$

Next,

$$E[U] = \int_0^1 \int_0^{1-u} 2u \, dv \, du \quad (13.0.162)$$

$$= \int_0^1 2u \int_0^{1-u} 1 \, dv \, du \quad (13.0.163)$$

$$= \int_0^1 2u(1-u) \, du \quad (13.0.164)$$

$$= 2 \int_0^1 (u-u^2) \, du \quad (13.0.165)$$

$$= 2 \left[\frac{1}{2}u^2 - \frac{1}{3}u^3 \right]_{u=0}^1 \quad (13.0.166)$$

$$= 2 \left(\frac{1}{2} - \frac{1}{3} \right) \quad (13.0.167)$$

$$= \frac{1}{3} \quad (13.0.168)$$

By symmetry, $E[V] = \frac{1}{3}$. Hence,

$$\text{Cov}[U, V] = E[UV] - E[U]E[V] \quad (13.0.169)$$

$$= \frac{1}{12} - \frac{1}{3} \frac{1}{3} \quad (13.0.170)$$

$$= \frac{1}{12} - \frac{1}{9} \quad (13.0.171)$$

$$= -\frac{1}{36}. \quad (13.0.172)$$

s.11.0.7. Since $(u, v) = g(x, y) = (x + y, x - y)$,

$$\frac{\partial(u, v)}{\partial(x, y)} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} = |-2| = 2.$$

Moreover, $x = (u + v)/2$, $y = (u - v)/2$, so that

$$f_{U,V}(u, v) = f_{X,Y}(x, y) \frac{\partial(x, y)}{\partial(u, v)} = f_{X,Y}(x, y)/2 = f_X(x)f_Y(y)/2 = \frac{1}{2}.$$

The difficulty is in the domain, however. Note that x and y satisfy $0 \leq x \leq 1$, $0 \leq y \leq 1$. So $0 \leq (u + v)/2 \leq 1$ and $0 \leq (u - v)/2 \leq 1$, which simplifies to $-v \leq u \leq 2 - v$ and $v \leq u \leq 2 + v$, which can also be written as $|v| \leq u \leq 2 - |v|$.

s.11.0.8. Note that $u(x, y) = \min\{x, y\} = xI_{x \leq y} + yI_{x > y}$. With a similar expression for v we find for the Jacobian:

$$\frac{\partial(u, v)}{\partial(x, y)} = \begin{bmatrix} I_{x \leq y} & I_{y < x} \\ I_{y < x} & I_{x \leq y} \end{bmatrix} = |I_{x \leq y} - I_{x > y}| = 1.$$

If $(U, V) = g(X, Y)$, then $g^{-1}(u, v) = \{(u, v), (v, u)\}$, i.e., a set of two points.

If X, Y iid with PDF, then $f_{X,Y}(x, y) = f(x)f(y)$.

s.11.0.9. a.

$$F \geq 0 \implies 2 < y \tag{13.0.173}$$

$$F \leq 1 \implies F(3, y) \leq 1 \implies F(3, 4) = 1 \tag{13.0.174}$$

b. $F(3, 7) = 1$.

c. $f(x, y) = \partial_x \partial_y F(x, y) = (x - 1)/4$ for $x \in (1, 3)$, $y \in (2, 4)$ and 0 elsewhere.

d.

$$P\{2 < X < 3\} = F_X(3) - F_X(2) \tag{13.0.175}$$

$$= F_{X,Y}(3, 4) - F_{X,Y}(2, 4) = 1 - 1 \cdot 2/8 = 3/4. \tag{13.0.176}$$

e. Make a drawing of the rectangle $[2, 3] \times [2, 4]$. Then check what parts of this are covered by $F_{X,Y}$.

$$P\{2 < X < 3, 2 < Y < 3\} = F_{X,Y}(3, 3) - F_{X,Y}(2, 3) - F_{X,Y}(3, 2) + F_{X,Y}(2, 2). \tag{13.0.177}$$

The rest is just number plugging.

f. Use the fundamental bridge and c.

$$P\{Y < 2X\} = E[I_{Y < 2X}] \quad (13.0.178)$$

$$= \iint I_{y < 2x} f_{X,Y}(x, y) dx dy \quad (13.0.179)$$

$$= \frac{1}{4} \iint I_{y < 2x} I_{2 < y < 4} I_{1 < x < 3} (x-1) dx dy \quad (13.0.180)$$

$$= \frac{1}{4} \int_1^3 (x-1) \int I_{2 < y < \min\{2x, 4\}} dy dx \quad (13.0.181)$$

$$= \frac{1}{4} \int_1^3 (x-1)(\min\{2x, 4\} - 2) dx \quad (13.0.182)$$

$$= \frac{1}{4} \int_1^2 (x-1)(2x-2) dx + \frac{1}{4} \int_2^3 (x-1)(4-2) dx. \quad (13.0.183)$$

Finishing the computation must be easy for you now (and if not, practice real hard).

g. As X, Y continuous, the answer is equal to that of f.

h. Similar to f. but a bit more involved.

$$P\{Y < 2X, Y + 2X > 6\} = E[I_{Y < 2X, Y > 6-2X}] \quad (13.0.184)$$

$$= \iint I_{y < 2x, y > 6-2x} f_{X,Y}(x, y) dx dy \quad (13.0.185)$$

$$= \frac{1}{4} \iint I_{y < 2x, y > 6-2x} I_{2 < y < 4} I_{1 < x < 3} (x-1) dx dy \quad (13.0.186)$$

$$= \frac{1}{4} \int_1^3 (x-1) \int I_{\max\{2, 6-2x\} < y < \min\{2x, 4\}} dy dx \quad (13.0.187)$$

$$= \frac{1}{4} \int_1^3 (x-1)[\min\{2x, 4\} - \max\{2, 6-2x\}]^+ dx, \quad (13.0.188)$$

where we need the $[\cdot]^+$ to ensure the positivity of $\min\{2x, 4\} - \max\{2, 6-2x\}$. To see this, make a graph of the function $\min\{2x, 4\} - \max\{2, 6-2x\}$. Also, from this graph,

$$= \frac{1}{4} \int_{3/2}^2 (x-1)(2x-6+2x) dx + \frac{1}{4} \int_2^3 (x-1)(4-2) dx. \quad (13.0.189)$$

The rest is for you.

s.11.0.10. The function $g(x) = x^4$ is not one-to-one on \mathbb{R} . It is, however, locally, one-to-one, around the roots of U . (In this course we don't deal with complex numbers, for your interest, it can be proven that the equation $x^4 - y$ has, in general, four roots in the complex plane.)

We need to be bit careful with applying the change of variables formula, but we are OK if we apply it locally around the roots $U^{1/4}$ and $-U^{1/4}$. However, mind that we also should take care of the domain of V , so it might be that these roots don't lie in the domain of V .

With all this, let's first tackle the Jacobian, and then get the domain right with indicators.

$$u = g(v) = v^4 \implies v = \pm u^{1/4}, \quad (13.0.190)$$

$$f_U(u)du = f_V(v)dv \implies f_U(u) = f_V(v) \frac{dv}{du}, \quad (13.0.191)$$

$$\frac{du}{dv} = 4v^3 = 4u^{3/4} I_{v \geq 0} - 4u^{3/4} I_{v < 0}, \quad (13.0.192)$$

$$f_U(u) = \frac{f_V(-u^{1/4})}{4(-u)^{3/4}} I_{-u^{1/4} \in (-3, 0)} + \frac{f_V(u^{1/4})}{4(u)^{3/4}} I_{u^{1/4} \in [0, 2)} \quad (13.0.193)$$

$$= \frac{f_V(-u^{1/4})}{4(-u)^{3/4}} I_{u \in (0, 81)} + \frac{f_V(u^{1/4})}{4(u)^{3/4}} I_{u \in [0, 16)}. \quad (13.0.194)$$

If V has the uniform distribution, then $f_V(v) = \frac{1}{5}$ for $v \in (-3, 2)$, so

$$f_U(u) = \frac{1}{20(-u)^{3/4}} I_{u \in (0, 81)} + \frac{1}{20(u)^{3/4}} I_{u \in [0, 16)}. \quad (13.0.195)$$

s.11.0.11. Here is a direct approach.

$$x = \tan u = g(u) \implies u = \arctan x \quad (13.0.196)$$

$$\frac{dx}{du} = \frac{1}{\cos^2 u} = \frac{\sin^2 u + \cos^2 u}{\cos^2 u} = \tan^2 u + 1 = x^2 + 1, \quad (13.0.197)$$

$$f_X(x) = f_U(u) \frac{du}{dx} = \frac{1}{\pi} I_{u \in (0, \pi)} \frac{1}{1+x^2} \quad (13.0.198)$$

$$= \frac{1}{\pi(1+x^2)} I_{\arctan x \in (0, \pi)} = \frac{1}{\pi(1+x^2)}. \quad (13.0.199)$$

In the last equation we just shifted the \tan from $(-\pi/2, \pi/2]$ to the interval $(0, \pi)$. The \tan has also a proper inverse in $(0, \pi)$ (make a drawing of \tan to see this), hence all is well-defined.

s.11.0.12. No. The only relevant information is the amount of legs won by each player.

s.11.0.13. Our current information can be represented as: $A_{10} = 6$.

s.11.0.14. We have $A_n \sim \text{Bin}(n, p)$.

s.11.0.15. Let f_0 denote the prior distribution of p . Then for the posterior pdf we find by Bayes' theorem:

$$f_1(p|A_n = k) = \frac{\mathbb{P}\{A_n = k | p\} f_0(p)}{\mathbb{P}\{A_n = k\}} \quad (13.0.200)$$

$$= \frac{\binom{n}{k} p^k (1-p)^{n-k} \cdot 1}{\mathbb{P}\{A_n = k\}} \quad (13.0.201)$$

$$\propto p^k (1-p)^{n-k}, \quad (13.0.202)$$

in which we recognize the pdf of a $\text{Beta}(k+1, n-k+1)$ distribution (up to a normalizing constant). Hence, $p|A_n = k \sim \text{Beta}(k+1, n-k+1)$.

s.11.0.16. Important: we have already observed 10 legs with an outcome, with which we have updated our belief. Hence, we should use the *posterior* distribution given $A_{10} = 6$ in this exercise! (It's easy to make a mistake here.) Think about it. Suppose instead we had observed 1000 legs and Amy had won 990 of them (i.e., $A_{1000} = 990$). Wouldn't we use this information if someone offered us a bet?

Note that Bob should win the match if and only if he wins the next three legs. Let W_k be short-hand notation for the event "Bob wins the k th leg". Then, observing that W_{11}, W_{12}, W_{13} are independent, and using the LOTP in the fourth step, we obtain

$$P\{\text{Bob wins the match} \mid A_{10} = 6\} = P\{W_{11} \cap W_{12} \cap W_{13} \mid A_{10} = 6\} \quad (13.0.203)$$

$$= P\{W_{11} \mid A_{10} = 6\} P\{W_{12} \mid A_{10} = 6\} P\{W_{13} \mid A_{10} = 6\} \quad (13.0.204)$$

$$= P\{W_{11} \mid A_{10} = 6\}^3 \quad (13.0.205)$$

$$= \int_0^1 P\{I_{11} \mid p, A_{10} = 6\}^3 f_1(p \mid A_{10} = 6) dp \quad (13.0.206)$$

$$= \int_0^1 (1-p)^3 \cdot \frac{p^6(1-p)^4}{\beta(7,5)} dp \quad (13.0.207)$$

$$= \frac{\beta(7,8)}{\beta(7,5)} \int_0^1 \frac{p^6(1-p)^7}{\beta(7,8)} dp \quad (13.0.208)$$

$$= \frac{\beta(7,8)}{\beta(7,5)}. \quad (13.0.209)$$

(Note that we very explicitly do all the steps here. It might be more intuitively clear if you skip the first few steps and write

$$P\{\text{Bob wins the match} \mid A_{10} = 6\} = \int_0^1 (1-p)^3 f_1(p \mid A_{10} = 6) dp, \quad (13.0.210)$$

and work from there).

s.11.0.17. We have

$$P\{\text{Bob wins the match} \mid A_{10} = 6\} = \frac{\beta(7,8)}{\beta(7,5)} \quad (13.0.211)$$

$$= \left(\frac{6!7!}{14!} \right) / \left(\frac{6!4!}{11!} \right) \quad (13.0.212)$$

$$= \frac{7!/4!}{14!/11!} \quad (13.0.213)$$

$$= \frac{7 \cdot 6 \cdot 5}{14 \cdot 13 \cdot 12} \quad (13.0.214)$$

$$= 5/52 \quad (13.0.215)$$

$$= 0.0962. \quad (13.0.216)$$

s.11.0.18. Our expected profit when taking the bet is

$$300 \cdot P\{\text{Bob wins the match} \mid A_{10} = 6\} - 10 \cdot P\{\text{Amy wins the match} \mid A_{10} = 6\} \quad (13.0.217)$$

$$= 300 \cdot \frac{5}{52} - 10 \cdot \left(1 - \frac{5}{52}\right) \quad (13.0.218)$$

$$= 19.808. \quad (13.0.219)$$

So we expect to make a profit of €19.81. Hence, you should take the bet.

s.11.0.19.

s.11.0.20. By the linearity of expectation and BH Theorem 9.3.9:

$$\begin{aligned} E[(Y - E[Y|X] - h(X))^2] &= E[(Y - E[Y|X])^2 - 2(Y - E[Y|X])h(X) + (h(X))^2] \\ &= E[(Y - E[Y|X])^2] - E[2(Y - E[Y|X])h(X)] + E[(h(X))^2] \\ &= E[(Y - E[Y|X])^2] + E[(h(X))^2]. \end{aligned}$$

Since $E[(h(X))^2] \geq 0$, we conclude that $E[(Y - E[Y|X] - h(X))^2] \geq E[(Y - E[Y|X])^2]$ for any function h , so $E[Y|X]$ is the predictor of Y based on X with the lowest mean squared error, i.e. the best predictor of Y based on X .

s.11.0.21.

$$E[\tilde{X} \mid Y] = E[X - \hat{X} \mid Y] = E[X \mid Y] - E[E[X \mid Y] \mid Y] \quad (13.0.220)$$

$$= E[X \mid Y] - E[X \mid Y] E[1 \mid Y] \quad (13.0.221)$$

$$= E[X \mid Y] - E[X \mid Y] 1 = 0 \quad (13.0.222)$$

s.11.0.22.

$$E[\tilde{X}] = E[E[\tilde{X} \mid Y]] = E[0 \mid Y] = 0. \quad (13.0.223)$$

This means that the estimation error \tilde{X} does not have bias.

s.11.0.23.

$$E[\tilde{X}\hat{X}] = E[E[\tilde{X}\hat{X} \mid Y]] \quad (13.0.224)$$

$$= E[E[\tilde{X} E[X \mid Y] \mid Y]] \quad (13.0.225)$$

$$= E[E[X \mid Y] E[\tilde{X} \mid Y]] \quad (13.0.226)$$

$$= E[E[X \mid Y] 0 \mid Y] = 0 \quad (13.0.227)$$

Here, in the rest of the exercises about this topic, we have seen the most terrible mistakes during grading. Hence, study the reasoning applied very carefully, and ensure you know the motivation behind each and every step. There will be questions in the exam about this, and you have to be able to use the arguments. If not, you fail the exam; simple as that. So, you are warned!

s.11.0.24. Using the previous exercises,

$$\text{Cov}[\hat{X}, \tilde{X}] = E[\hat{X}\tilde{X}] - E[\hat{X}]E[\tilde{X}] = 0 - E[\hat{X}]0 = 0. \quad (13.0.228)$$

Next, from the definition of $\tilde{X} = X - \hat{X} \implies X = \tilde{X} + \hat{X}$. The variance of the sum is the sum of the variances since \hat{X} and \tilde{X} are uncorrelated.

s.11.0.25. Since $E[\tilde{X}] = 0$,

$$V[\tilde{X}] = E[\tilde{X}^2] \quad (13.0.229)$$

$$= E[E[\tilde{X}^2 | Y]] \quad (13.0.230)$$

$$= E[E[(X - \hat{X})^2 | Y]] \quad (13.0.231)$$

$$= E[E[(X - E[X|Y])^2 | Y]] \quad (13.0.232)$$

$$= E[V[X|Y]], \quad (13.0.233)$$

where the last equation follow from the definition of $V[X|Y]$.

For Eve's law, use the above and the previous exercise to see that

$$V[X] = V[\hat{X}] + V[\tilde{X}] = V[E[X|Y]] + E[V[X|Y]]. \quad (13.0.234)$$

s.11.0.26. From Probability Theory we know $E\left(\frac{Y_n}{n}\right) = \frac{1}{2}$ and $V\left(\frac{Y_n}{n}\right) = \frac{1}{4n}$. Then by Chebyshev's inequality,

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{Y_n}{n} - \frac{1}{2}\right| \geq \varepsilon\right) \leq \lim_{n \rightarrow \infty} \frac{V\left(\frac{Y_n}{n}\right)}{\varepsilon^2} = \lim_{n \rightarrow \infty} \frac{1}{4n\varepsilon^2} = 0 \quad \forall \varepsilon > 0.$$

s.11.0.27. The Cauchy distribution has no mean to converge to.