

# Homework 2: Probability (2)

UA CSC 380: Principles of Data Science, Spring 2023

Homework due at 11:59pm on Feb 8

**Deliverables** You must make two submissions: (1) your homework as a SINGLE PDF file by the stated deadline to the gradescope. Include your code and output of the code as texts in the PDF. and (2) your codes as a ZIP file to a separate submission. Each subproblem is worth 10 points. More instructions:

- You can hand-write your answers and scan them to make it a PDF. If you use your phone camera, I recommend using TurboScan (smartphone app) or similar ones to avoid uploading a slanted image or showing the background. Make sure you rotate it correctly.
- Watch the video and follow the instruction for the submission: [https://youtu.be/KMPoby5g\\_nE](https://youtu.be/KMPoby5g_nE)
- **Show all work along with answers to get the full credit.**
- Place your final answer into an ‘answer box’ that can be easily identified.
- There will be no late days. Late homeworks result in zero credit.

Failure to follow the submission instruction will result in a minor penalty in credit.

**You can choose to work individually or in pairs.**

- If you choose to work in pairs, you are free to discuss whatever you want with your partner; please make only one submission per group.
- Please do not discuss with people outside your group about the homework (refer to the academic integrity policy in Lecture 1).
- If you have clarification questions, please feel free to post on Piazza so that it can promote discussion.

## Problem 1: Joint, Conditional, Marginal Probability

Suppose we throw a fair six-sided die twice in a row. Let  $A$  be a random variable representing the number on the first throw, and  $B$  be the number on the second throw. Let  $M = A \times B$  be the product of both throws. What are the following probabilities?

a)  $P(M = 5)$

b)  $P(M = 4)$

c)  $P(A = 2, B = 3 \mid M = 6)$

d)  $P(M = 6 \mid A = 2, B = 3)$

e)  $P(A = 1, B = 3 \mid M = 6)$

f)  $P(A = 4 \mid M = 6)$

g)  $P(A = 1 \mid M = 6)$

h)  $P(A = 3 \mid M = 9)$

## Problem 2: Discrete Approximation

In continuous probability, we often need to solve messy integrals. For example, in this class we might need to use integrals to evaluate the probability of an event under a cumulative distribution function (CDF). Rather than solve this by hand, we can approximate it using discrete intervals. This problem will explore discrete approximation of integrals using a Gaussian model. Recall that the probability density function of a Gaussian is,

$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right).$$

In the questions below, we will use Python to form a discrete approximation of this continuous distribution, and evaluate associated probabilities.

- a) Form a discrete approximation of the Normal PDF with mean  $\mu = 70$  and standard deviation  $\sigma = 2$ . To do this, create an array  $x$  of evenly spaced values in the range  $[68, 76]$  at increments of 2 excluding 76 (this array will include 68 and 74). The function `numpy.arange` might be helpful. Create an array  $p$  containing values of the PDF at each location  $x$ . Plot the result as a bar chart (use `matplotlib.pyplot.bar`). In the same figure, overlay a PDF curve (use `matplotlib.pyplot.plot`) at more finely spaced intervals (e.g. 0.01). Paste your code here.
- b) The bar chart above is a discrete approximation of the continuous PDF. We will use it to approximate  $P(68 < X \leq 76)$ . Recall that  $\mathcal{N}(x | \mu, \sigma^2)$  is the PDF of  $X$ , so

$$P(68 < X \leq 76) = \int_{68}^{76} \mathcal{N}(x | \mu, \sigma^2) dx.$$

We will approximate this integral using a Riemann sum ([https://en.wikipedia.org/wiki/Riemann\\_sum](https://en.wikipedia.org/wiki/Riemann_sum)). Let  $N$  be the number of grid points in your array  $x$ . The spacing between grid points is  $\Delta x$  and let the  $i^{\text{th}}$  point of array  $p$  be  $p_i$ . The Riemann sum approximation is,

$$P(68 < X \leq 76) \approx \sum_{i=1}^N p_i \Delta x$$

What is the value of the Riemann sum approximation to  $P(68 < X \leq 76)$ ? Paste your code here.

- c) Now, reduce the spacing  $\Delta x = 0.01$  and recompute the discrete approximation of  $P(68 < X \leq 76)$ . How do the two approximations compare? What is the practical downside of smaller spacing?
- d) Repeat the steps above to show the distribution over the range  $[20, 120]$  and compute  $P(20 \leq X < 120)$ . What is the value? This interval should contain almost all of the probability in this distribution, i.e. the event is almost certain.

### Problem 3: Conditional Independence

Let  $A, B, C \in \{0, 1\}$  be three binary random variables with the following joint probability distribution:

$a$	$b$	$c$	$P(A = a, B = b, C = c)$
0	0	0	0.01
0	0	1	0.07
0	1	0	0.02
0	1	1	0.10
1	0	0	0.02
1	0	1	0.30
1	1	0	0.04
1	1	1	0.44

- By direct calculation, compute the marginal  $P(A, B)$  (recall that  $P(A, B)$  is represented by 4 numbers:  $P(A = 0, B = 0)$ ,  $P(A = 0, B = 1)$ ,  $P(A = 1, B = 0)$ ,  $P(A = 1, B = 1)$ ).
- By direct calculation compute the marginals  $P(A)$  and  $P(B)$ .
- Are the random variables  $A$  and  $B$  independent? Why or why not?
- Compute the conditional distribution  $P(A, B | C)$ . Note that this includes computing  $P(A, B | C = 0)$  as well as  $P(A, B | C = 1)$ , each of which is represented by 4 numbers (in total 8 numbers).
- Calculate  $P(A = 0, B = 0 | C = 1)$  as well as  $P(A = 0 | C = 1) \cdot P(B = 0 | C = 1)$ . Are  $A$  and  $B$  conditionally independent given  $C$ ? Why or why not?

## Problem 4: Diagnostic Tests and Bayes' Rule

I have decided to get myself tested for COVID-19 antibodies. However, being comfortable with statistics, I am curious about what the test means for my actual status. Let's investigate these questions, showing all your work.

- a) The antibody test I take has a *sensitivity* (a.k.a. true positive rate) of 97.5% and a *specificity* (a.k.a. true negative rate) of 99.1%. If you are not familiar with sensitivity vs specificity, please see Wikipedia. Assume that 4% of the population actually have COVID-19 antibodies. Write down the joint probability distribution  $P(S, R)$  with events for antibody state  $S \in \{\text{true}, \text{false}\}$  and test result  $R \in \{\text{true}, \text{false}\}$ .
- b) Assuming I receive a *positive* test result, use Bayes' rule to calculate the probability that I actually have COVID-19 antibodies.
- c) Assuming I receive a *negative* test result, what is the probability that I *do not* have COVID-19 antibodies?
- d) Assume I take the test twice, and receive a positive result in the first test and a negative result in the second test. Assume that the two test results are conditionally independent given the existence of the antibody. What is the probability that I have COVID-19 antibodies according to Bayes' rule?
- e) Now assume that only 1% of the population has COVID-19 antibodies. Repeat parts (b) and (c) with this revised prior belief.