

## *Foundations of Statistics*

### **Solutions to Homework 2**

#### **Topic: Conditional probability and independence**

##### **Part I. theoretical problems**

In this section, let  $(\Omega, \mathcal{A}, \mathbb{P})$  be a probability space.

**1.** Suppose  $A, B \in \mathcal{A}$  are events with  $0 < \mathbb{P}(A) < 1$  and  $0 < \mathbb{P}(B) < 1$ .

**(a)** If  $A$  and  $B$  are disjoint, can they be independent?

*Solution:* No.  $A$  and  $B$  are disjoint. Thus  $A \cap B = \emptyset$  and

$$0 = \mathbb{P}(A \cap B) \neq \mathbb{P}(A)\mathbb{P}(B) > 0.$$

**(b)** If  $A$  and  $B$  are independent, can they be disjoint?

*Solution:* No. We have

$$\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B) > 0,$$

therefore,  $A \cap B \neq \emptyset$ , i.e.,  $A$  and  $B$  are not disjoint.

**(c)** If  $A \subset B$ , can  $A$  and  $B$  be independent?

*Solution:* No. We have

$$\mathbb{P}(A \cap B) = \mathbb{P}(A) > \mathbb{P}(A)\mathbb{P}(B),$$

as  $1 > \mathbb{P}(B)$ . Therefore,  $A$  and  $B$  cannot be independent.

**(d)** If  $A$  and  $B$  are independent, can  $A$  and  $A \cup B$  be independent?

*Solution:* No. We suppose  $A$  and  $A \cup B$  are independent and show that this leads to a contradiction.

$$\begin{aligned}\mathbb{P}(A) &= \mathbb{P}(A \cap (A \cup B)) \\ &= \mathbb{P}(A)\mathbb{P}(A \cup B) \implies \mathbb{P}(A \cup B) = 1\end{aligned}$$

On the other hand, we have

$$\begin{aligned}\mathbb{P}(A \cup B) &= \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B) \\ &= \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A)\mathbb{P}(B) \\ &= \mathbb{P}(A)(1 - \mathbb{P}(B)) + \mathbb{P}(B)\end{aligned}$$

Taking these two into account, we conclude

$$1 - \mathbb{P}(B) = \mathbb{P}(A)(1 - \mathbb{P}(B)) \implies \mathbb{P}(A) = 1$$

which is in contradiction with the original assumption  $\mathbb{P}(A) < 1$ . So  $A$  and  $A \cup B$  can not be independent.

**2.** Let  $B \in \mathcal{A}$  be an event with  $\mathbb{P}(B) > 0$ . Prove that  $\mathbb{Q}(\cdot) := \mathbb{P}(\cdot|B)$  is a probability measure on  $(\Omega, \mathcal{A})$ . In other words, show that

**(a)**  $\mathbb{Q}(\Omega) = 1$ .

*Solution:* We have

$$\mathbb{Q}(\Omega) = \mathbb{P}(\Omega|B) = \frac{\mathbb{P}(\Omega \cap B)}{\mathbb{P}(B)} = \frac{\mathbb{P}(B)}{\mathbb{P}(B)} = 1.$$

**(b)** For any countable family of mutually disjoint sets  $(A_n)_{n=1}^{\infty}$  with  $A_n \in \mathcal{A}$ , we have  $\mathbb{Q}(\cup_{n=1}^{\infty} A_n) = \sum_{n=1}^{\infty} \mathbb{Q}(A_n)$ . This means

$$\mathbb{P}\left(\bigcup_{n=1}^{\infty} A_n \mid B\right) = \sum_{n=1}^{\infty} \mathbb{P}(A_n \mid B).$$

*Solution:* We have

$$\begin{aligned}
\mathbb{P}\left(\bigcup_{n=1}^{\infty} A_n | B\right) &= \frac{\mathbb{P}\left(\left(\bigcup_{n=1}^{\infty} A_n\right) \cap B\right)}{\mathbb{P}(B)} \\
&= \frac{\mathbb{P}\left(\bigcup_{n=1}^{\infty} (A_n \cap B)\right)}{\mathbb{P}(B)} \\
&\stackrel{*}{=} \frac{\sum_{n=1}^{\infty} \mathbb{P}(A_n \cap B)}{\mathbb{P}(B)} \\
&= \sum_{n=1}^{\infty} \frac{\mathbb{P}(A_n \cap B)}{\mathbb{P}(B)} \\
&= \sum_{n=1}^{\infty} \mathbb{P}(A_n | B)
\end{aligned}$$

where in the step  $*$ , we used the fact that sets are disjoint. More precisely,  $A_n \cap A_m = \emptyset$  for  $m \neq n$  also implies

$$(A_n \cap B) \cap (A_m \cap B) = A_n \cap A_m \cap B = \emptyset.$$

**3.** Suppose  $A, B \in \mathcal{A}$  are events with  $\mathbb{P}(B) > 0$ .

(a) Use exercise 2 to conclude that  $\mathbb{P}(A|B) + \mathbb{P}(A^c|B) = 1$ .

*Solution:*  $A$  and  $A^c$  are obviously disjoint. By applying the results of exercise 2, we obtain

$$\begin{aligned}
\mathbb{P}(A|B) + \mathbb{P}(A^c|B) &\stackrel{2(b)}{=} \mathbb{P}(A \cup A^c | B) \\
&= \mathbb{P}(\Omega | B) \stackrel{2(a)}{=} 1.
\end{aligned}$$

(b) Give counterexamples to show that in general the following statements are false:

(i)  $\mathbb{P}(A|B) + \mathbb{P}(A|B^c) = 1$ ,

*Solution:* Take  $A = \Omega$

$$\mathbb{P}(A|B) + \mathbb{P}(A|B^c) = \mathbb{P}(\Omega|B) + \mathbb{P}(\Omega|B^c) = 1 + 1 = 2.$$

(ii)  $\mathbb{P}(A|B) + \mathbb{P}(A^c|B^c) = 1$ .

*Solution:* Take  $A = B$

$$\mathbb{P}(A|B) + \mathbb{P}(A^c|B^c) = \mathbb{P}(B|B) + \mathbb{P}(B^c|B^c) = 1 + 1 = 2.$$

4. Let  $0 < \mathbb{P}(B) < 1$  and  $\mathbb{P}(A|B^c) = \mathbb{P}(A|B)$ .  
Show that the events  $A$  and  $B$  must be independent.

*Solution:* We have

$$\begin{aligned}\frac{\mathbb{P}(A \cap B^c)}{\mathbb{P}(B^c)} &= \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} \\ \frac{\mathbb{P}(A) - \mathbb{P}(A \cap B)}{1 - \mathbb{P}(B)} &= \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} \\ \mathbb{P}(A)\mathbb{P}(B) &= \mathbb{P}(A \cap B)\end{aligned}$$

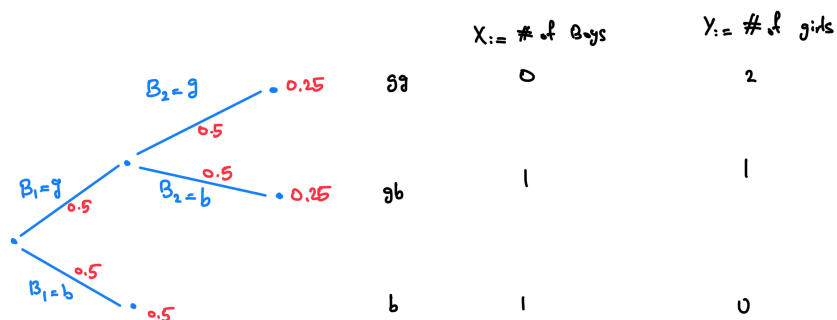
which means that the events  $A$  and  $B$  are independent.

## Part II. practical problems

**5 (demographic problem).** The “one child rule” in some provincial parts of China had been changed to the following. All couples are allowed one baby. If the baby is girl, they are allowed to have exactly one more. If this rule is exactly followed (and ignoring possibilities of twins, etc.), what will be the resulting proportion of boys to girls in this community? Assume that for each child the probability of being girl is 0.5.

*Hint:* draw a tree diagram to determine the corresponding probabilities and then calculated the expected value for boys and girls in a typical family.

*Solution:*



$$\begin{aligned}
 \left( \mathbb{P}(B_1 = g) = 0.5 \right) \times \left( \mathbb{P}(B_2 = g | B_1 = g) = 0.5 \right) &= \left( \mathbb{P}(B_1 = g \cap B_2 = g) \right) = 0.25 \\
 \left( \mathbb{P}(B_1 = g) = 0.5 \right) \times \left( \mathbb{P}(B_2 = b | B_1 = g) = 0.5 \right) &= \left( \mathbb{P}(B_1 = g \cap B_2 = b) \right) = 0.25 \\
 \left( \mathbb{P}(B_1 = b) = 0.5 \right) &= 0.5
 \end{aligned}$$

In addition, we have

$$\begin{aligned}
 \mathbb{E}[X] &= 0.25 \times 0 + 0.25 \times 1 + 0.5 \times 1 = 0.75 \\
 \mathbb{E}[Y] &= 0.25 \times 2 + 0.25 \times 1 + 0.5 \times 0 = 0.75
 \end{aligned}$$

So the resulting proportion of boys to girls in this community will stay 1:1.

**6. (epidemiologic problem).** For the Roche Sars-CoV-2 Antigen Rapid Test, the following information is provided by the manufacturer on its accuracy:

- **Sensitivity** = 96.52%
- **Specificity** = 99.68%.

**Sensitivity** is the conditional probability of a positive test when there is infection with Covid-19, and **specificity** is the conditional probability of a negative test, provided there is no infection with Covid.

In the following we consider the events

- “+” = {positive test}, “**C**” = {Covid Infection}.
- “−” = {negative test}, “**N**” = {No Covid Infection}.

*Solution:* We have

$$\begin{aligned}\mathbb{P}(+|C) &= 0.9652 && \text{(sensitivity)} \\ \mathbb{P}(-|N) &= 0.9968 && \text{(specificity)}\end{aligned}$$

		actual status	
		C	N
test result	+	true positive	false positive
	-	false negative	true negative

**(a)** What is the probability of a (false) positive test in a non-infected person?

*Solution:* We have

$$\mathbb{P}(+|N) = 1 - \mathbb{P}((+)^c|N) = 1 - \mathbb{P}(-|N) = 1 - 0.9968 = 0.0032.$$

**(b)** Calculate the probability that a randomly selected person will test positive, given that 1 in 150 people in your area are currently infected with Covid (which is close to the actual 7-day incidence in Bielefeld and NRW on 21.10.2022).

*Solution:* What we know and what we want to find is the following:

$$\mathbb{P}(C) = \frac{1}{150} \quad \implies \quad \mathbb{P}(+) = ?$$

To this end, let us first recall the *law of total probability*:

Given a probability space  $(\Omega, \mathcal{A}, \mathbb{P})$ , if  $(A_n)_{n=1}^N \subset \mathcal{A}$  are disjoint sets with  $\cup_{n=1}^N A_n = \Omega$ , then for any arbitrary set  $B \in \mathcal{A}$ , we have

$$\mathbb{P}(B) = \sum_n \mathbb{P}(B|A_n)\mathbb{P}(A_n).$$

Therefore, we have

$$\begin{aligned}\mathbb{P}(+) &= \mathbb{P}(+|C)\mathbb{P}(C) + \mathbb{P}(+|N)\mathbb{P}(N) \\ &= 0.9652 \times \frac{1}{150} + 0.0032 \times \frac{149}{150} \\ &= 0.009613\end{aligned}$$

(c) Now use Bayes' theorem to calculate the probability that a person who tests positive really has Covid (under the assumptions made in (b)).

*Solution:* We have

$$\begin{aligned}\mathbb{P}(C|+) &= \frac{\mathbb{P}(+|C)\mathbb{P}(C)}{\mathbb{P}(+)} \quad (\text{Bayes' theorem}) \\ &= \frac{0.9652 \times \frac{1}{150}}{0.009613} = 66.94\%\end{aligned}$$

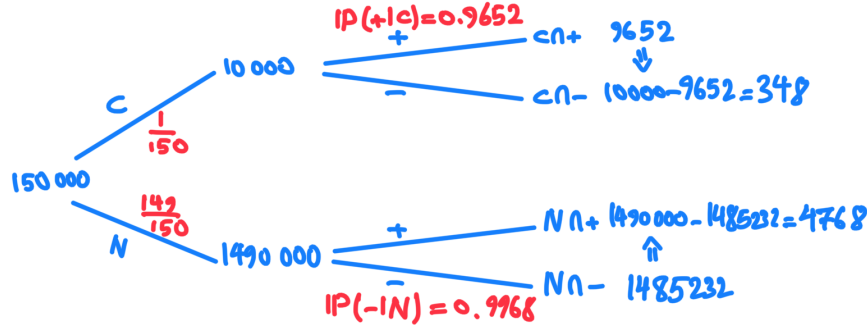
(d) Conversely, if a person tests negative, what is the probability of being infected anyway?

*Solution:* We have

$$\begin{aligned}\mathbb{P}(C|-) &= \frac{\mathbb{P}(-|C)\mathbb{P}(C)}{\mathbb{P}(-)} \quad (\text{Bayes' theorem}) \\ &= \frac{(1 - 0.9652) \times \frac{1}{150}}{(1 - 0.009613)} = 0.0234\%\end{aligned}$$

(e) Check your result from (c) using a probability tree in which you assume that 1 500 000 randomly selected persons ( $\rightsquigarrow$  how many of them have Covid/no Covid  $\rightsquigarrow$  how many of these will test positive/negative in turn).

*Solution:*



Based on tree diagram above, we can check part (c):

$$\# \text{ people with } + = 9652 + 4768 = 14420$$

$$\# \text{ people with } + \text{ and } C = 9652$$

therefore,

$$\mathbb{P}(C|+) = \frac{\mathbb{P}(+ \cap C)}{\mathbb{P}(+)} = \frac{9652}{14420} = 66.94\%$$

(f) Plot the diagnostic power of the test (i.e.  $\mathbb{P}(C|+)$ ) as a function of the actual incidence (i.e.  $\mathbb{P}(C)$ ) in R (both axis in %) and describe the dependence. Mark the point corresponding to the probability calculated in part (c) with **BI2020**. By repeating the same calculations, find the power of test if 1 in 10 people were infected with Covid (this e.g. corresponds to the incidence rate in New York in Spring 2020). Mark the point again in the plot with **NY2020**.

*Solution:* By combining the formulas of part (b) and (c), we obtain

$$\begin{aligned} \mathbb{P}(C|+) &\stackrel{(c)}{=} \frac{\mathbb{P}(+|C)\mathbb{P}(C)}{\mathbb{P}(+)} \\ &\stackrel{(b)}{=} \frac{\mathbb{P}(+|C)\mathbb{P}(C)}{\mathbb{P}(+|C)\mathbb{P}(C) + \mathbb{P}(+|N)\mathbb{P}(N)} \\ &= \frac{\mathbb{P}(+|C)\mathbb{P}(C)}{\mathbb{P}(+|C)\mathbb{P}(C) + (1 - \mathbb{P}(-|N))(1 - \mathbb{P}(C))} \end{aligned}$$

So we have found how  $\mathbb{P}(C|+)$  depends on the actual incidence  $\mathbb{P}(C)$ . We emphasize that this function only depends on the parameters  $\mathbb{P}(+|C)$  (sensitivity) and  $\mathbb{P}(-|N)$  (specificity), which are determined by the manufacturer. In



particular,

$$\mathbf{BI2020} : \quad \mathbb{P}(C) = \frac{1}{150} \quad \implies \quad \mathbb{P}(C|+) = 66.94\%$$

$$\mathbf{NY2020} : \quad \mathbb{P}(C) = \frac{1}{10} \quad \implies \quad \mathbb{P}(C|+) = 97.10\%$$

(g) At what proportion of actually infected persons would a positive test be equivalent to a 50/50 situation, i.e.  $\mathbb{P}(C|+) = 0.5$ ? Add a horizontal line into your plot and mark the corresponding point with a star.

*Solution:* Here  $\mathbb{P}(C|+) = 50\%$ . Solving the Eq. above we get  $P(C) \approx \frac{1}{302}$ .

(h) Confirm your numerical results in (c) by doing computer simulation with R (see page 26 of Ch. 1.2 in the lecture notes).

*Solution:* See next page.

```

# Exercise 6
# part (f)
pC <- 0.9652 # sensitivity
nN <- 0.9968 # specificity

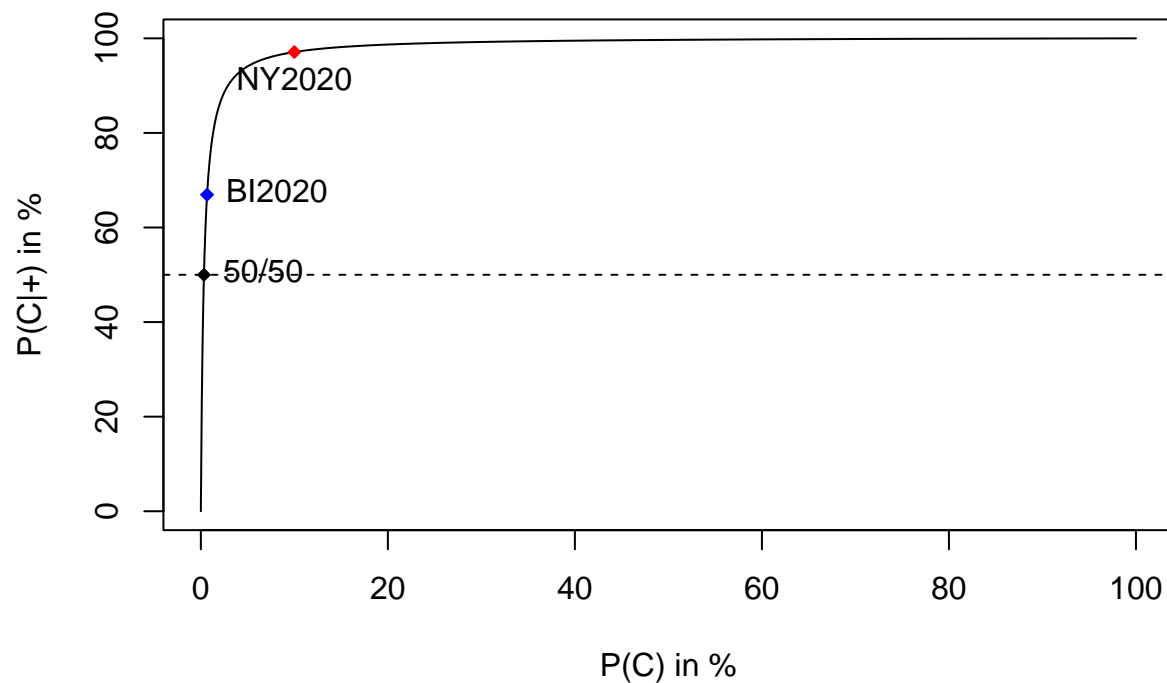
C <- seq(0,1,by=0.001)
Cp <- (pC * C)/(pC * C + (1-nN) * (1 - C))
plot(C*100,Cp*100,type="l",xlab="P(C) in %",ylab = "P(C|+) in %")

# mark the points
C_BI <- 1/150
Cp_BI <- (pC * C_BI)/(pC * C_BI + (1-nN) * (1 - C_BI))
points(C_BI*100, Cp_BI*100, col = "blue", pch = 18)
text(C_BI*100, Cp_BI*100, labels = "BI2020", pos = 4)

C_NY <- 1/10
Cp_NY <- (pC * C_NY)/(pC * C_NY + (1-nN) * (1 - C_NY))
points(C_NY*100, Cp_NY*100, col = "red", pch = 18)
text(C_NY*100, Cp_NY*100, labels = "NY2020", pos = 1)

# part (g)
abline(h = 50, col = "black", lty = 2)
C_50 <- 1/302.62 # by solving the eq.
Cp_50 <- (pC * C_50)/(pC * C_50 + (1-nN) * (1 - C_50))
points(C_50*100, Cp_50*100, col = "black", pch = 18)
text(C_50*100, Cp_50*100, labels = "50/50", pos = 4)

```



```

# part (h)
nloop=100000
numerator=0; denominator=0

c <- 1/150 # also try 1/10

for (iloop in 1:nloop){
  disease=sample(0:1,1,prob=c(1-c, c)) ## sample a person
  if(disease==0) {
    test=sample (0:1,1,prob=c(0.9968, 1-0.9968)) ## if no disease
  }else{
    test=sample(0:1,1,prob=c(1-0.9652, 0.9652)) ## if disease
  }
  if(test==1){
    denominator=denominator+1
    if(disease==1){ numerator=numerator+1}
  }
}
numerator/denominator * 100 # computes  $P(C|+) = P(C \text{ and }+)/P(+)$ 

## [1] 66.42857

```