

## Homework 2

Professor: Ziyu Shao

Due: 2022/09/25 10:59pm

1. **(Bertrands Box Paradox)** There are three boxes:

- (a) A box containing two gold coins
- (b) A box containing two silver coins
- (c) A box containing one gold coin and a silver coin

After choosing a box at random and withdrawing one coin at random, if that happens to be a gold coin, find the probability of the next coin drawn from the same box also being a gold coin.

2. Alice is trying to communicate with Bob, by sending a message (encoded in binary) across a channel.

- (a) Suppose for this part that she sends only one bit (a 0 or 1), with equal probabilities. If she sends a 0, there is a 5% chance of an error occurring, resulting in Bob receiving a 1; if she sends a 1, there is a 10% chance of an error occurring, resulting in Bob receiving a 0. Given that Bob receives a 1, what is the probability that Alice actually sent a 1?
- (b) To reduce the chance of miscommunication, Alice and Bob decide to use a repetition code. Again Alice wants to convey a 0 or a 1, but this time she repeats it two more times, so that she sends 000 to convey 0 and 111 to convey 1. Bob will decode the message by going with what the majority of the bits were. Assume that the error probabilities are as in (a), with error events for different bits independent of each other. Given that Bob receives 110, what is the probability that Alice intended to convey a 1?

3. To battle against spam, Bob installs two anti-spam programs. An email arrives, which is either legitimate (event  $L$ ) or spam (event  $L^c$ ), and which program  $j$  marks as legitimate (event  $M_j$ ) or marks as spam (event  $M_j^c$ ) for  $j \in \{1, 2\}$ . Assume that 10% of Bobs email is legitimate and that the two programs are each “90% accurate” in the sense that  $P(M_j|L) = P(M_j^c|L^c) = 9/10$ . Also assume that given whether an email is spam, the two programs’ outputs are conditionally independent.

- (a) Find the probability that the email is legitimate, given that the 1st program marks it as legitimate (simplify).

- (b) Find the probability that the email is legitimate, given that both programs mark it as legitimate (simplify).
- (c) Bob runs the 1st program and  $M_1$  occurs. He updates his probabilities and then runs the 2nd program. Let  $\tilde{P}(A) = P(A|M_1)$  be the updated probability function after running the 1st program. Explain briefly in words whether or not  $\tilde{P}(L|M_2) = P(L|M_1 \cap M_2)$ : is conditioning on  $M_1 \cap M_2$  in one step equivalent to first conditioning on  $M_1$ , then updating probabilities, and then conditioning on  $M_2$ ?
4. (a) Suppose that in the population of college applicants, being good at baseball is independent of having a good math score on a certain standardized test (with respect to some measure of “good”). A certain college has a simple admissions procedure: admit an applicant if and only if the applicant is good at baseball or has a good math score on the test.
- Give an intuitive explanation of why it makes sense that among students that the college admits, having a good math score is negatively associated with being good at baseball, i.e., conditioning on having a good math score decreases the chance of being good at baseball.
- (b) Show that if  $A$  and  $B$  are independent and  $C = A \cup B$ , then  $A$  and  $B$  are conditionally dependent given  $C$  (as long as  $P(A \cap B) > 0$  and  $P(A \cup B) < 1$ ), with

$$P(A | B, C) < P(A | C).$$

This phenomenon is known as Berkson’s paradox, especially in the context of admissions to a school, hospital, etc.

5. We want to design a spam filter for email. A major strategy is to find phrases that are much more likely to appear in a spam email than in a no spam email. In that exercise, we only consider one such phrase: “free money”. More realistically, suppose that we have created a list of 100 words or phrases that are much more likely to be used in spam than in non-spam. Let  $W_j$  be the event that an email contains the  $j$ th word or phrase on the list. Let

$$p = P(\text{spam}), p_j = P(W_j|\text{spam}), r_j = P(W_j|\text{not spam})$$

where “spam” is shorthand for the event that the email is spam.

Assume that  $W_1, \dots, W_{100}$  are conditionally independent given that the email is spam, and also conditionally independent given that it is not spam. A method for classifying emails (or other objects) based on this kind of assumption is called a *naive Bayes classifier*. (Here “naive” refers to the fact that the conditional independence is a strong assumption, not to Bayes being naive. The assumption may or may not be realistic, but naive Bayes classifiers sometimes work well in practice even if the assumption is

not realistic.)

Under this assumption we know, for example, that

$$P(W_1, W_2, W_3^c, W_4^c, \dots, W_{100}^c | spam) = p_1 p_2 (1 - p_3)(1 - p_4) \dots (1 - p_{100}).$$

Without the naive Bayes assumption, there would be vastly more statistical and computational difficulties since we would need to consider  $2^{100} \approx 1.3 \times 10^{30}$  events of the form  $A_1 \cap A_2 \dots \cap A_{100}$  with each  $A_j$  equal to either  $W_j$  or  $W_j^c$ . A new email has just arrived, and it include the 23rd, 64th, and 65th words or phrases on the list (but not the other 97). So we want to compute

$$P(spam | W_1^c, \dots, W_{22}^c, W_{23}, W_{24}^c, \dots, W_{63}^c, W_{64}, W_{65}, W_{66}^c, \dots, W_{100}^c).$$

Note that we need to condition on *all* the evidence, not just the fact that  $W_{23} \cap W_{64} \cap W_{65}$  occurred. Find the condition probability that the new email is spam (in terms of  $p$  and the  $p_j$  and  $r_j$ ).