

CS2402 Assignment 2

Lecturer: Kede Ma & Shuaicheng Li

Due on Apr 7, 2024 23:59

You should submit your solution to canvas before the deadline. Late submission halves the score.

*For all the problems in this assignment, you **SHOULD** present your solution step by step instead of giving the answers only.*

Question 1 (10pts)

Independent random variables X and Y follow distributions $\mathcal{N}(4,11)$ and $\mathcal{N}(6,14)$, respectively. Show that the probability that $(X - Y)^2$ has a value between 4 and 16 is approximately 0.25.

Question 2 (10pts)

In a sequence of Bernoulli trials with probability p of success, find the probability that a successes will occur before the b^{th} failures. (Suppose $a = 2$, $b = 3$, s represents success and f represents failures, then the events of 2 successes will occur before the 3rd failures include: $ssfff$, $sfsff$, $sffsf$, $fssff$, $fsfsf$, $ffssf$.)

Question 3 (10pts)

- (a) The random variable X is defined as the larger value obtained in two throws of an unbiased six-sided die. Show that the probability mass function is

$$p(x) = \frac{2x-1}{36}, x = 1, 2, \dots, 6.$$

- (b) The random variable Y represents the highest value obtained in k throws of an unbiased six-sided die. Find an expression for the probability mass function of Y .

[Hint: find the distribution first.]

Question 4(10pts)

Consider a small town where each family can have 0, 1 or 2 children. The probability of gender is equal and independent. Let the probability of a family of zero or one child be r , and the probability of having a second child be $1 - r$. Calculate the following probabilities:

- (a) $P(\text{a family has two girls})$.
- (b) $P(\text{the elder child is a boy and the younger child is a girl})$.
- (c) $P(\text{a family has at least one boy})$.

Please show the problem solving process and answers.

Question 5 (10pts)

Suppose a diagnostic test for a rare disease has the following characteristics:

- The probability that a person has the disease is 1% (prior probability).
- The probability that the test correctly identifies a sick person as having the disease is 99%.
- The probability that the test correctly identifies a healthy person as not having the disease is 95%.

- (a) Fill in the table below with the probability (expressed as a percentage) of each situation:

	Disease Present (D+)	Disease Absent (D-)
Test Positive (T+)	TP: ____	FP: ____
Test Negative (T-)	FN: ____	TN: ____

- (b) Calculate the posterior probability that a patient actually has the disease when given the positive test result.

Question 6 (10pts)

Suppose you have two fair coins and you flip them both. Let X be the number of heads that show. Then, you grab a standard deck of 52 cards and draw X cards with replacement. Let Y be the number of aces you draw.

- (a) Find $P(Y = 1|X = 2)$.
- (b) Find $P(Y = 0)$, $P(Y = 1)$, and $P(Y = 2)$.
- (c) Find $P(X = 2|Y = 1)$.

Question 7 (10pts)

A producer of a certain type of electronic component ships to suppliers in lots, each lot has twenty components. Suppose that 60% of all such lots contain no defective components, 30% contain one defective component, and 10% contain two defective components. A lot is picked, two components from the lot are randomly selected and tested, and neither is defective.

- (a) What is the probability that zero defective components exist in the lot?
- (b) What is the probability that one defective exists in the lot?
- (c) What is the probability that two defectives exist in the lot?

Question 8 (10pts)

The table below shows some data from the early days of a clothing company. Each row shows the sales for a year, and the amount spent on advertising in that year.

Year	Advertising (Million Dollars)	Sales (Million Dollars)
1	22	650
2	31	855
3	33	1064
4	44	1191
5	51	1420

- (a) Use linear regression to calculate the relationship between the advertising and sales (slope, and intercept).
- (b) Calculate the correlation coefficient between the amounts of advertising and sales.
- (c) If there is no advertising, what is the expected sales? If the amount of advertising is 58 million dollars in the next year, please predict the sales.

Question 9 (10pts)

Let $\theta > 0$, and consider the p.d.f. $p(x) = \theta e^{-\theta x}, x \geq 0$. Based on a random sample X_1, \dots, X_n from this distribution, find the maximum likelihood estimator for θ .

Hint: $a \ln \theta - b\theta$ achieves its maximum when $\theta = \frac{a}{b}$.

Question 10 (10pts) Implementation

You need to implement a python program detector.py (included in assignment2 fold) to solve the problem. However, you only need to include the solutions with blanks filled out in your submitted PDF file.

Use Bayes' rule to detect the spam email. The prior probability of spam emails are 80%.

The key words are listed in the following table

Word	P(word spam)	P(word \neg spam)
\$	88%	47%
donate	66%	10%
research	11%	60%
contact	51%	51%
CS2402	0.5%	40%

The two emails are:

Email I.

Dear research student, we have just uploaded the competition questions for CS2402. Please sign up and contact me before this Friday (29, March, 2024) if you have any questions. The award for the winner is \$10,000,000! Also, you can donate your award to the whole class.

Email II.

I have decided to donate what I have to you. I was diagnosed with cancer of the lungs few years ago. I have been inspired by God to donate my inheritance to you for the good work of God and charity purpose.

I am doing this because my family are unbelievers and I will not allow them inherit this money for their own selfishness. I decided to bequeath the sum of \$10,000,000.00 to you. If you are much more interested, Contact Thomas with this specified email: thmasbfd@gmail.com).

Python sketch codes for detecting spam email is provided in “detector.py”. You need to implement Bayes rule to compute the probability of spam email for the given two emails.

- (a) Please fill in the blanks and execute “detector.py” with the two exemplified emails, respectively.
- (b) Please complete the tables.

detector.py is included in assignment2 fold

Hint: copy the email and paste it to the shell

```
email = input('Please enter the test Email: ')    # gets the test email
email = email.lower()                           #convert string to lowercase
print(email)
```

Input the spam detection keyword

```
word1 = input('Please enter the first key word: ') # gets the first keyword
word1 = word1.lower()                            #convert the first key word to lowercase

word2 = _____                              # gets the second keyword
word2 = word2.lower()                            # convert the second keyword to lowercase
word3 = _____                              # gets the third keyword
word3 = word3.lower()                            # convert the third keyword to lowercase

word4 = _____                              # gets the fourth keyword
word4 = word4.lower()                            # convert the fourth keyword to lowercase

word5 = _____                              # gets the fifth keyword
```

```

word5 = word5.lower()                # convert the fifth keyword to lowercase

# Detect whether the key word occurred in your email #

n_word1 = email.count(word1) # count the number of the first keyword in the email
n_word2 = email.count(word2) # count the number of the second keyword in the email
n_word3 = email.count(word3) # count the number of the third keyword in the email
n_word4 = email.count(word4) # count the number of the fourth keyword in the email
n_word5 = email.count(word5) # count the number of the fifth keyword in the email

# The prior odds of spam is 80% #

p_spam = 1.0                        # initialize the likelihood ratio of the spam email
p_no_spam = 1.0                    # initialize the likelihood ratio of the no spam email
prior_odds_spam = _____ # prior odds of spam

# Calculate the Likelihood ratio #

if n_word1 != 0:                    # the first keyword occurred in the email
    p_spam = p_spam * 0.88
    p_no_spam = p_no_spam * 0.47

if n_word2 != 0:                    # the second keyword occurred in the email
    p_spam = _____
    p_no_spam = _____

if n_word3 != 0:                    # the third keyword occurred in the email
    p_spam = _____
    p_no_spam = _____

if n_word4 != 0:                    # the fourth keyword occurred in the email
    p_spam = _____
    p_no_spam = _____

```

```

if n_word5 != 0:                # the fifth keyword occurred in the email

    p_spam = _____

    p_no_spam = _____

# Calculate the Posterior odds #

posterior_odds_spam = _____    # posterior odds of spam

posterior_odds_no_spam = _____    # posterior odds of no spam


p_isSpam = _____    # probability of spam P(spam | email)


print("p_spam: %.6f, p_no_spam: %.6f, posterior_odds_spam: %.6f, posterior_odds_no_spam: %.6f,
p_isSpam: %.6f"%(p_spam, p_no_spam, posterior_odds_spam, posterior_odds_no_spam, p_isSpam))

```

Email I:

Hypothesis	Prior odds	Likelihood ratio	Posterior odds	P(spam email)
spam	0.8			
¬spam	0.2			

Email II:

Hypothesis	Prior odds	Likelihood ratio	Posterior odds	P(spam email)
spam	0.8			
¬spam	0.2			