

This problem set covers material from Week 10, dates 4/23- 4/25. Unless otherwise noted, all problems are taken from the textbook. Problems can be found at the end of the corresponding subsection.

Instructions: Write or type complete solutions to the following problems and submit answers to the corresponding Canvas assignment. Your solutions should be neatly-written, show all work and computations, include figures or graphs where appropriate, and include some written explanation of your method or process (enough that I can understand your reasoning without having to guess or make assumptions). A general rubric for homework problems appears on the final page of this assignment.

Feel free to reference previous results where appropriate!

Tuesday 4/23

1. Section 9.1: problem 10 (note you did problem 8 on previous homework.)
2. Suppose $X|\theta \sim \text{Unif}[0, \theta]$. We'd like to test the following hypotheses:

$$H_0 : \theta \leq 1 \quad \text{vs.} \quad H_1 : \theta > 1$$

Let δ_c be the procedure that rejects H_0 if $X \geq c$ for some $c > 0$. For each possible value of $X = x$, find the form of the p-value if $X = x$ is observed.

Thursday 4/24

3. Suppose X_1, \dots, X_n are iid $N(\mu, 1)$ with $\mu \in \mathbb{R}$ unknown. We are interested in testing the following hypotheses:

$$H_0 : \mu = 0 \quad H_1 : \mu \neq 0$$

In the exercise, you will construct a likelihood ratio test for these hypotheses.

- (a) Recall that the likelihood ratio statistic is defined as

$$\Lambda(\mathbf{X}) = \frac{\sup_{\mu \in \Omega_0} f(\mathbf{x}|\mu)}{\sup_{\mu \in \Omega} f(\mathbf{x}|\mu)}$$

Find formulas for the numerator and the denominator of the likelihood ratio statistic.

- (b) Show that the likelihood ratio statistic $\Lambda(\mathbf{x})$ can be simplified as

$$\Lambda(\mathbf{x}) = \exp\left(-\frac{n(\bar{x})^2}{2}\right)$$

Hint: Use one of our favorite identities.

- (c) Recall that a likelihood ratio test is of the form: “Reject H_0 if $\Lambda(\mathbf{x}) \leq k$,” where k is a constant with $0 \leq k \leq 1$. Show that the rejection region for a likelihood ratio test can also be expressed as

$$R = \left\{ |\bar{X}| \geq \sqrt{\frac{-2 \log k}{n}} \right\}$$

In other words, a likelihood ratio test is of the form “Reject H_0 if \bar{X} is far from 0.”

- (d) Find a value of k so that the likelihood ratio test has size $\alpha = 0.05$.
4. Let $X_1, \dots, X_n | \theta \stackrel{\text{iid}}{\sim} \text{Exp}(\theta)$. Suppose we have the following hypotheses:

$$H_0 : \theta = \theta_0 \quad \text{vs.} \quad H_1 : \theta \neq \theta_0$$

- (a) Derive the likelihood ratio test of the above hypotheses. Then show that the rejection region is of the form $\{\bar{X} e^{-\theta_0 \bar{X}} \leq c\}$ for some $c > 0$.
- (b) For the remainder of this problem suppose that $\theta_0 = 1$, $n = 15$, and we want a level-0.05 test. In order to use this test, we must find the appropriate value of c .

Argue (visually and/or using calculus and/or mathematical reasoning) that an equivalent rejection region is of the form $\{\bar{X} \leq c_0\} \cup \{\bar{X} \geq c_1\}$ where c_0 and c_1 are determined by c .

- (c) Explain why c should be chosen so that $\Pr(\bar{X} e^{-\bar{X}} \leq c) = 0.05$.
- (d) Recall that since the X_i are iid Exponential, \bar{X} is itself a gamma distributed random variable. Using this knowledge, explain as clearly as possible how you could approximate the value of c . *Hint: consider using some sort of simulation.*
- (e) Copy-and-paste the following into R:

```
x <- c(0.843, 0.100, 1.025, 0.018, 0.143, 0.315, 1.554, 0.091, 0.922,
0.585, 0.155, 0.087, 0.141, 0.275, 1.567)
```

Based on this observed data \mathbf{x} , the testing procedure outlined above, and your method to find c , determine whether or not we should reject H_0 .

Friday 4/25

5. Hardy-Weinberg revisited. Recall the two-allele scenario: an individual can have any of the two following genotypes: AA , Aa , or aa . If we randomly select n individuals from a population, we count up how many of the individuals have each allele. Thus we observe $(N_{AA}, N_{Aa}, N_{aa}) \sim \text{Multinom}(n, (p_{AA}, p_{Aa}, p_{aa}))$. Clearly, the three probabilities must be non-zero and sum to 1.

Under the assumption of Hardy-Weinberg (H-W) equilibrium, allele A is present in the population with unknown probability θ , and allele a with probability $1 - \theta$, such that

$$p_{AA} = \theta^2 \quad p_{Aa} = 2\theta(1 - \theta) \quad p_{aa} = (1 - \theta)^2$$

In practice, we do not know if H-W is present in a given population, but we might use data to test the following hypotheses:

$$H_0 : p_{AA} = \theta^2, p_{Aa} = 2\theta(1 - \theta), p_{aa} = (1 - \theta)^2 \quad (\text{i.e. H-W is present})$$

$$H_1 : p_{AA}, p_{Aa}, p_{aa} \text{ are any three non-zero probabilities such that } p_{AA} + p_{Aa} + p_{aa} = 1 \\ (\text{i.e. H-W not present})$$

We will test these two hypotheses using a likelihood ratio test.

- Find formulas for the numerator and the denominator of the likelihood ratio statistic.
- Obtain the likelihood ratio statistic for these hypotheses. Simplify as much as possible.
- Recall the following data from Homework 3: In a sample from the Chinese population of Hong Kong in 1937, blood types occurred with the following frequencies, where A and a are erythrocyte antigens:

AA	Aa	aa	Total
342	500	187	1029

On Homework 3, many of you said you believed H-W was present in this population based on some MLEs. Let's test this more rigorously: what is the p -value of the observed data for this testing procedure? *Hint: note that since n is large and the Multinomial distribution satisfies "regularity conditions", you may be able to use an asymptotic result...*

- Based on your findings in (c), re-answer the following question: do you believe Hardy-Weinberg is present in this population? Briefly explain why or why not.
6. Suppose $X_1, \dots, X_n | \mu \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$ where μ is unknown and σ^2 is known. Consider the following hypotheses:

$$H_0 : \mu = \mu_0 \quad \text{vs.} \quad H_1 : \mu \neq \mu_0$$

Let the test statistic be $r(\mathbf{X}) = |\bar{X} - \mu_0|$, and the procedure δ defines the rejection region $R = [c, \infty)$ (i.e. reject H_0 if $r(\mathbf{X}) \geq c$).

- If we want δ to be a level- α_0 test, what should c be? *Good practice for midterm!*
 - Use your answer from (a) to obtain a $(1 - \alpha_0)$ coefficient confidence interval for μ .
7. Prove the following theorem:
- Duality Thm. 2:** Let $c(\mathbf{X})$ be a $\gamma = 1 - \alpha_0$ coefficient confidence set for θ . Define the acceptance region R^c for a test δ_{θ_0} of $H_0 : \theta = \theta_0$ as $R^c = \{r(\mathbf{x}) : \theta_0 \in c(\mathbf{X})\}$. Then δ_{θ_0} is level- α_0 .

General rubric

Points	Criteria
5	The solution is correct <i>and</i> well-written. The author leaves no doubt as to why the solution is valid.
4.5	The solution is well-written, and is correct except for some minor arithmetic or calculation mistake.
4	The solution is technically correct, but author has omitted some key justification for why the solution is valid. Alternatively, the solution is well-written, but is missing a small, but essential component.
3	The solution is well-written, but either overlooks a significant component of the problem or makes a significant mistake. Alternatively, in a multi-part problem, a majority of the solutions are correct and well-written, but one part is missing or is significantly incorrect.
2	The solution is either correct but not adequately written, or it is adequately written but overlooks a significant component of the problem or makes a significant mistake.
1	The solution is rudimentary, but contains some relevant ideas. Alternatively, the solution briefly indicates the correct answer, but provides no further justification.
0	Either the solution is missing entirely, or the author makes no non-trivial progress toward a solution (i.e. just writes the statement of the problem and/or restates given information).
Notes:	For problems with multiple parts, the score represents a holistic review of the entire problem. Additionally, half-points may be used if the solution falls between two point values above.
Notes:	For problems with code, well-written means only having lines of code that are necessary to solving the problem, as well as presenting the solution for the reader to easily see. It might also be worth adding comments to your code.