

This problem set covers material from Week 2, dates 2/16- 2/20. 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.

Unless otherwise stated, you must confirm that your critical point is indeed a maximum for full credit!

Tuesday 2/16

- Let $X_1, \dots, X_n | \theta \stackrel{\text{iid}}{\sim} \text{Poisson}(\theta)$. In class, we showed $\hat{\theta}_{MLE} = \bar{X}$ given X_1, \dots, X_n . Now define the following random variable:

$$Y_i = \begin{cases} 0 & \text{if } X_i = 0 \\ 1 & \text{if } X_i > 0 \end{cases}$$

We seek to find an MLE for θ using the Y_1, \dots, Y_n .

- Find the PMF for Y_i and the likelihood function for θ given y_1, y_2, \dots, y_n . (*Hint:* it should have a similar structure to the Bernoulli likelihood function.)

- Show that

$$\hat{\theta}_{MLE} = \log \left(\frac{n}{\sum(1 - y_i)} \right)$$

is the MLE estimate for θ given y_1, \dots, y_n and provided that not all the y_i 's are equal to 1.

- What does the likelihood function equal if $y_i = 1$ for all $i = 1, \dots, n$? Use the likelihood to explain why the MLE for θ does not exist if all of the y_i 's are equal to 1.
- For fixed n , what the probability that all of the y_i 's are equal to 1? What happens to that probability as $\theta \rightarrow \infty$? What does that say about the probability that MLE does not exist as $\theta \rightarrow \infty$?
- For a fixed θ , what happens to the probability that the MLE does not exist as $n \rightarrow \infty$?

- Hardy-Weinberg equilibrium.** The Hardy-Weinberg equilibrium is a principle stating that the genetic variation in a population will remain constant from one generation to the next in the absence of disturbing factors. In the simplest case, there are two kinds of alleles denoted A and a with frequencies θ and $1 - \theta$, respectively. Organisms

inherit two alleles (one from the father and one from the mother), and the possible genotypes (pairs of alleles) an offspring can have are AA , Aa , and aa . The different ways to form these genotypes from a male and female parent are as follows, where the quantity in the brackets [] denotes the proportion/frequency *under the assumption of Hardy-Weinberg equilibrium*:

		Female	
		A [θ]	a [$1 - \theta$]
Male	A [θ]	AA [θ^2]	Aa [$\theta(1 - \theta)$]
	a [$1 - \theta$]	Aa [$\theta(1 - \theta)$]	aa [$(1 - \theta)^2$]

Note that the genotype frequencies sum to one: $\theta^2 + 2\theta(1 - \theta) + (1 - \theta)^2 = 1$.

In a random sample of n independent observations, let n_{AA} denote the number of observations with the AA genotype, with analogous definitions for n_{Aa} and n_{aa} , and so $n = n_{AA} + n_{Aa} + n_{aa}$.

- (a) Under the assumption of Hardy-Weinberg equilibrium (i.e. that AA occurs with probability θ^2 , Aa with probability $2\theta(1 - \theta)$ and aa with probability $(1 - \theta)^2$), what is the probability $P(N_{AA} = n_{AA}, N_{Aa} = n_{Aa}, N_{aa} = n_{aa})$?
 - (b) Under the assumption of Hardy-Weinberg equilibrium, obtain an MLE for θ .
3. We will work with a non-continuous/discrete parameter space in this problem!

The capture/recapture method in ecology is a way to estimate the unknown size of an animal population in the wild. We will denote the unknown population size as parameter N . This is how capture/recapture works:

- Researchers go out and capture M individuals, mark them, and release them back into the wild.
- After an appropriate period of time, they perform a second capture event where they capture K individuals. Some of these may be re-captures of the M previously marked individuals.
- We'll denote X as the number of re-captured marked individuals among the K .

Not super important, but the method assumes that each animal in the population is equally likely to be captured, and no individuals move in/out of the study region or give birth/die during the study period.

We get to observe data (M, K, X) , and we'd like to use them to obtain an MLE for N .

- (a) Let's consider X (the number of individuals that were tagged when re-captured) the random variable of interest. Formulate and justify a Hypergeometric statistical model for X in terms of M, K, N .

- (b) To find an MLE in this case, we want to find the value of N that maximizes $\mathcal{L}(N|x, M, K)$. We can't go through our usual calculus methods because the likelihood is discrete! Instead, we can find the value of largest value of N such that

$$\frac{\mathcal{L}(N|x, M, K)}{\mathcal{L}(N-1|x, M, K)} \geq 1$$

(assuming the likelihood is concave). Briefly explain why this is an appropriate way to find the MLE.

- (c) We could find \hat{N}_{MLE} using the method identified in (b), but that's ugly and involves lots of factorials! Instead, let's obtain the maximum likelihood estimate empirically using a plot for a specific set of data. Suppose $M = 12$ animals are initially captured, tagged and released. A few months later, $K = 20$ animals are captured, examined, and released. $X = 6$ of these were tagged. Using R, make a plot of both the likelihood function and the log-likelihood function for N . Using your plot/code, what is \hat{N}_{MLE} ?

Wednesday 2/18

4. Let $X_1, \dots, X_n | \boldsymbol{\theta} \stackrel{iid}{\sim} f(x|\alpha, \beta)$ where $\boldsymbol{\theta} = (\alpha, \beta)$ with corresponding CDF

$$F(x|\alpha, \beta) = \Pr(X_i \leq x|\alpha, \beta) = \begin{cases} 0 & \text{if } x < 0 \\ \left(\frac{x}{\beta}\right)^\alpha & \text{if } 0 \leq x \leq \beta \\ 1 & \text{if } x > \beta, \end{cases}$$

where the parameters α and β are positive and unknown. Find the MLE of $\boldsymbol{\theta} = (\alpha, \beta)$.

The direct way of using the Hessian won't be useful here (think about why!). Instead, obtain the MLE by doing the following:

- For fixed α , what value of β maximizes the likelihood? Call this value $\hat{\beta}(\alpha)$ (it may or may not depend on α).
 - Then, find the value of α that maximizes the likelihood when we fix $\beta = \hat{\beta}(\alpha)$.
5. 7.5: Problem 12 *Note: the implied statistical model is a Multinomial. The sum-to-one constraint on the proportions is important! Dig deep into your Calculus classes and try to re-learn the method of LaGrange Multipliers. Also, you do not need to verify that the critical point is indeed a maximum; believe me that it is!*

Friday 2/20

6. 7.6: Problem 4 (*Hint: what is the statistical/sampling model for X , and how does that relate to T and/or β ?*)

7. 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

Suppose we are interested in the true proportions of each of these three blood types among the Chinese population in Hong Kong at that time. That is, we would like to estimate the true proportions of AA , Aa , and aa blood types.

- (a) Without assuming Hardy-Weinberg equilibrium, obtain maximum likelihood estimates of the three proportions from these data. Provide brief justification of how you obtained your MLEs, if applicable.
 - (b) Assuming Hardy-Weinberg equilibrium, obtain maximum likelihood estimates of the three proportions from these data. Provide brief justification of how you obtained your MLEs, if applicable.
 - (c) How do the sets of estimates in (a) and (b) compare? Do you believe that there is Hardy-Weinberg equilibrium for these blood types? Why or why not?
8. Suppose $X_1, \dots, X_n | \theta \stackrel{\text{iid}}{\sim} \text{Unif}[0, \theta]$. Show that the sequence of MLEs for θ , $\hat{\theta}_n = \max\{X_1, \dots, X_n\}$ is a consistent sequence of estimators for θ .

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).
<hr/>	
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.