

# Data 102, Spring 2023 Midterm 1

- You have 110 minutes to complete this exam. There are 5 questions, totaling 40 points.
- You may use one  $8.5 \times 11$  sheet of handwritten notes (front and back), and the provided reference sheet. No other notes or resources are allowed.
- You should write your solutions inside this exam sheet.
- You should write your name and Student ID on every sheet (in the provided blanks).
- Make sure to write clearly. We can't give you credit if we can't read your solutions.
- Even if you are unsure about your answer, it is better to write down partial solutions so we can give you partial credit.
- We have provided two blank pages of scratch paper, one at the beginning and one at the end of the exam. No work on these pages will be graded.
- You may, without proof, use theorems and facts that were given in the discussions or lectures, **but please cite them.**
- We don't answer questions individually. If you believe something is unclear, bring your question to us and if we find your question valid we will make a note to the whole class.
- Unless otherwise stated, no work or explanations will be graded for multiple-choice questions.
- Unless otherwise stated, you must show your work for free-response questions in order to receive credit.

Last name	
First name	
Student ID (SID) number	
Berkeley email	
Name of person to your left	
Name of person to your right	

## *Honor Code*

I will respect my classmates and the integrity of this exam by following this honor code.

I affirm:

- All of the work submitted here is my original work.
- I did not collaborate with anyone else on this exam.

Signature: \_\_\_\_\_

1. (5 points) For each of the following, determine whether the statement is true or false. For this question, no work will be graded and no partial credit will be assigned.
- (a) (1 point) The choice of the constant  $M$  or proposal distribution  $f(x)$  in rejection sampling has no effect on sampling efficiency, as long as  $Mq(x) \leq f(x)$ , where  $q(x)$  is the unnormalized target density.
- A. TRUE    B. FALSE
- (b) (1 point) Consider two medical tests for a disease. The first test has TPR=0.9 and FPR=0.05, while the second test has TPR=0.54 and FPR=0.03. Then the first test will always have a higher FDR than the second test.
- A. TRUE    B. FALSE
- (c) (1 point) Given specific sample data, the Benjamini-Hochberg procedure guarantees that the FDP will be lower than the requested level  $\alpha$ .
- A. TRUE    B. FALSE
- (d) (1 point) When using a GLM to fit continuous  $y$  with the Bernoulli likelihood, the Identity and Sigmoid are valid choices for the inverse link function.
- A. TRUE    B. FALSE
- (e) (1 point) When conducting linear regression in Bayesian perspective, the choice of prior will determine the form of regularization applied to the model.
- A. TRUE    B. FALSE

2. (8 points) **A different approach to FWER control.** Consider the following algorithm, known as the Holm-Bonferroni procedure:

1. Given a significance level  $\alpha \in [0, 1]$  and a set of  $n$  p-values,  $p_1, \dots, p_n$ . Sort the p-values in non-decreasing order:  $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(n)}$
  2. For  $k \in \{1, 2, \dots, n\}$ , if  $p_{(k)} \leq \frac{\alpha}{n-k+1}$ , reject the corresponding null hypothesis and continue. Otherwise, fail to reject all remaining hypotheses.
- (a) (2 points) **For this part only**, we consider the following 5 p-values for multiple hypothesis testing:

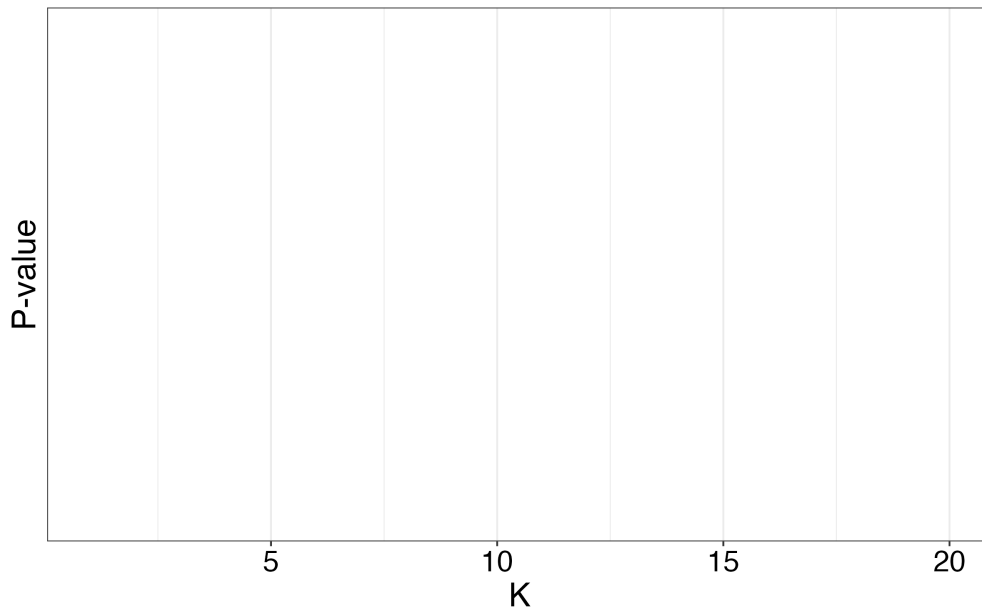
p-value	threshold	decision	reality
0.001			1
0.007			1
0.01			0
0.1			0
0.16			0

Fill in the threshold and decision columns of the above table for the Holm-Bonferroni procedure with level  $\alpha = 0.05$ . How many tests does the procedure reject?

- (b) (1 point) Like the Bonferroni correction, the Holm-Bonferroni procedure controls the family-wise error rate at level  $\alpha$ . Does the Holm-Bonferroni method make more or less discoveries than the Bonferroni correction? Justify your answer.

(c) (1 point) Comparing the Benjamini-Hochberg procedure with Bonferroni, which one makes fewer discoveries? In other words, does Bonferroni at rate  $\alpha$  also controls FDR at rate  $\alpha$  or Benjamini-Hochberg at rate  $\alpha$  also controls FWER at rate  $\alpha$ ? Justify your answer in words.

(d) (2 points) Assuming  $n = 20$ , draw the Benjamini-Hochberg guide line and Holm-Bonferroni guide line ( $\frac{\alpha}{n-k+1}$ ) on the same plot. X-axis should be  $k$  and the y-axis should be the p-value threshold. The Holm-Bonferroni guideline does not need to be exact, but its shape and position relative to the BH line should be accurate. Make sure you specify the equation for each line.

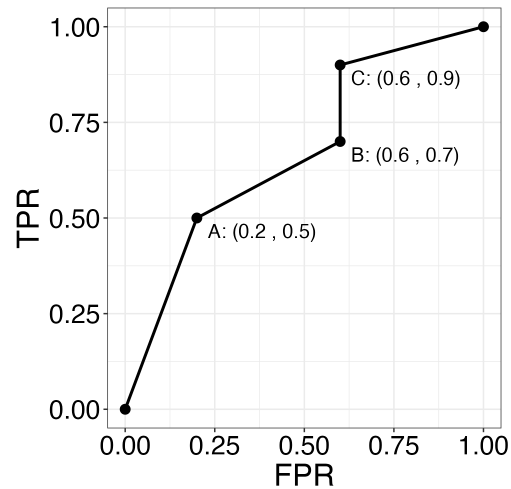


Name: \_\_\_\_\_

SID: \_\_\_\_\_

- (e) (2 points) Comparing the Benjamini-Hochberg procedure with Holm-Bonferroni, which one makes more discoveries for the same significance level  $\alpha$ ? You must show your work: show *mathematically* that either all discoveries made by Benjamini-Hochberg will also be made by Holm-Bonferroni, or the opposite.

3. (9 points) Your friend has developed a new cancer detection algorithm based on imaging and plans to evaluate its performance using an ROC curve. After testing the algorithm on samples from many patients, your friend generates the following ROC curve with the important points labeled with their corresponding (FPR, TPR). Throughout, assume that a positive case corresponds to having cancer.

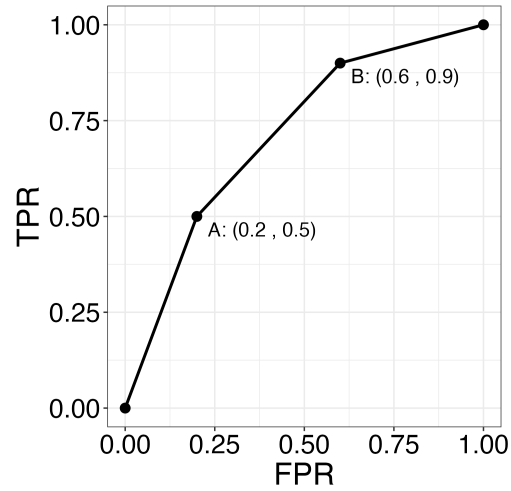


- (a) (1 point) What are the FNR and TNR associated with point B?

- (b) (2 points) Fill in the blanks below and explain your answer.

Among the points A, B, and C, \_\_\_\_\_ is strictly better than \_\_\_\_\_.

- (c) (3 points) It is possible to modify the algorithm and obtain the following ROC curve instead. First, explain why the modified algorithm is better: in other words, explain what measure we can use to make such a comparison. Second, describe how we can obtain the improved ROC curve using the algorithm which generated the ROC curve in parts (a), (b).



- (d) (3 points) A hospital looks into using this algorithm, and determines the cost of incorrectly classifying a cancer patient as not having cancer is \$1000, whereas the cost of incorrectly classifying a non-cancer patient as having cancer \$100. What should the baseline prevalence of cancer be such that you are indifferent between points A and B in the modified ROC curve from part (c)?



4. (10 points) Joe has landed a summer internship job in a customer service department. His job is to model the number of complaints received per week. He obtains data corresponding to  $n$  weeks  $\{x_1, x_2, \dots, x_n\}$  where  $x_i$  is the number of complaints received in week  $i$ . Each  $x_i$  follows a Poisson distribution with parameter  $\lambda$ :

$$x_i \sim \text{Poisson}(\lambda)$$

The PMF of a Poisson random variable with parameter  $\lambda$  is:  $P(x = k) = \frac{\lambda^k e^{-\lambda}}{k!}$ . You can assume that the number of complaints received in different weeks are independent of each other:  $x_i \perp\!\!\!\perp x_j \forall i \neq j$ .

- (a) (2 points) Suppose he wants to conduct the following hypothesis test:

$$H_0 : \lambda = \lambda_1$$

$$H_1 : \lambda = \lambda_2$$

where  $\lambda_2 > \lambda_1$ . He needs to fix his significance level at  $\alpha$ . Find the decision rule of the most powerful test for his problem. Your decision rule should fill in the blank in the following sentence with a mathematical expression that depends on  $x_1, \dots, x_n, \alpha, \lambda_1, \lambda_2$  and other constants: "If \_\_\_\_\_, then reject the null hypothesis".

*Hint: you don't need to simplify your expression or solve for the exact rejection threshold: you can just express it as a function that depends on the constant(s). Make sure you specify which constant(s) affects the threshold.*

(b) (2 points) Derive the Maximum Likelihood Estimator for  $\lambda$ .

(c) (1 point) Joe asks for your help to set up the problem from a Bayesian perspective. He makes the following choices:

- $\lambda$  is distributed according to a Gamma distribution with parameters  $\alpha > 0$  and  $\beta > 0$ :  $\lambda \sim \text{Gamma}(\alpha, \beta)$ .
- Each  $x_i$  is still distributed according to a Poisson with parameter  $\lambda$  and is conditionally independent of other  $x_j$ 's given  $\lambda$ .

Draw the graphical model for the setup above.

- (d) (2 points) Using the Bayesian model in the previous part, derive the posterior distribution for  $\lambda$  after observing  $\{x_1, \dots, x_n\}$ . The PDF of a Gamma distribution with parameters  $\alpha$  and  $\beta$  has the following form:

$$p(y; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} y^{\alpha-1} e^{-\beta y}$$

where  $\Gamma(\alpha)$  is a function that depends on  $\alpha$  which is a constant here. The reference sheet includes more information on the Gamma distribution.

*Hint: you should work with the unnormalized posterior which is proportional to the true posterior. You don't need to carry over the constants.*

- (e) (2 points) Using the posterior distribution from the previous part, find the MAP and MMSE estimate of  $\lambda$ .

- (f) (1 point) Compare the MMSE estimate from previous part with the MLE from above. Are they identical? If not, when would they become identical?

5. (8 points) Consider the following model with unknown variables  $\lambda$  and  $\mu$ , observed variables  $x_1, \dots, x_n$  and known constants  $\alpha$  and  $\beta$

$$\begin{aligned}\lambda &\sim \text{Gamma}(\alpha, \beta) \\ \mu \mid \lambda &\sim \text{Exponential}(\lambda) \\ x_i \mid \mu, \lambda &\sim \mathcal{N}(\mu, \lambda^2)\end{aligned}$$

- (a) (1 point) For this part only, suppose  $n = 2$ . Draw a graphical model for the variables described above.

- (b) (3 points) The following pseudocode provides a description of Gibbs sampling, but it contains exactly two mistakes. Circle each mistake, and in the space below, write the correct version (if the correction is just to remove the circled part, just write “remove only”)

*Hint: the fixes involve removing or changing part of the algorithm: no moving around is necessary.*

- (a) Compute the distributions  $p(\lambda \mid \mu)$  and  $p(\mu \mid \lambda, x_1, \dots, x_n)$
- (b) Initialize the following variables to zero:  $\lambda, \mu, x_1, \dots, x_n$
- (c) Repeat the following steps until enough samples have been obtained:
  - (i) Using the current values of  $\lambda$  and  $x_1, \dots, x_n$ , draw a sample for  $\mu$  from the conditional distribution in step (a).
  - (ii) Using the current values of  $\mu$  and  $x_1, \dots, x_n$ , draw a sample for  $\lambda$  from the other conditional distribution in step (a).
  - (iii) Save the current values of  $\lambda$  and  $\mu$  as samples.

- (c) (2 points) In step c-(ii) of the Gibbs sampling procedure above, we need to obtain a new sample for  $\lambda$ . Choose the single most efficient sampling algorithm to use to approximate the distribution for  $\lambda$ , or if sampling is not necessary, select option C. Ensure to justify your answer.
- A. Rejection sampling
  - B. Metropolis-Hastings
  - C. Sampling is not necessary

- (d) (2 points) Suppose we notice we have made a mistake in our model and the variance of each  $x$  is actually a known constant  $\sigma^2$ . In other words, our new model is:

$$\begin{aligned}\lambda &\sim \text{Gamma}(\alpha, \beta) \\ \mu \mid \lambda &\sim \text{Exponential}(\lambda) \\ x_i \mid \mu &\sim \mathcal{N}(\mu, \sigma^2)\end{aligned}$$

with unknown variables  $\lambda$  and  $\mu$ , observed variables  $x_1, \dots, x_n$  and known constants  $\alpha$ ,  $\beta$  and  $\sigma^2$ .

Now under this new model, what is the most efficient algorithm for sampling  $\lambda$  in step c-(ii) of the Gibbs sampling procedure above? Ensure to justify your answer.

- A. Rejection sampling
- B. Metropolis-Hastings
- C. Sampling is not necessary

# Midterm 1 Reference Sheet

## Algorithm 1 The Benjamini-Hochberg Procedure

**input:** FDR level  $\alpha$ , set of  $n$  p-values  $P_1, \dots, P_n$

Sort the p-values  $P_1, \dots, P_n$  in non-decreasing order  $P_{(1)} \leq P_{(2)} \leq \dots \leq P_{(n)}$

Find  $K = \max\{i \in \{1, \dots, n\} : P_{(i)} \leq \frac{\alpha}{n}i\}$

Reject the null hypotheses (declare discoveries) corresponding to  $P_{(1)}, \dots, P_{(K)}$

### Useful Distributions:

Distribution	Support	PDF/PMF	Mean	Variance	Mode
$X \sim \text{Poisson}(\lambda)$	$k = 0, 1, 2, \dots$	$\frac{\lambda^k e^{-\lambda}}{k!}$	$\lambda$	$\lambda$	$[\lambda]$
$X \sim \text{Gamma}(\alpha, \beta)$	$x \geq 0$	$\frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}$	$\frac{\alpha}{\beta}$	$\frac{\alpha}{\beta^2}$	$\frac{\alpha-1}{\beta}$
$X \sim \mathcal{N}(\mu, \sigma^2)$	$x \in \mathbb{R}$	$\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right)$	$\mu$	$\sigma^2$	$\mu$
$X \sim \text{Exponential}(\lambda)$	$x \geq 0$	$\lambda \exp(-\lambda x)$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$	0

**Conjugate Priors:** For observations  $x_i, i = 1, \dots, n$ :

Likelihood	Prior	Posterior
$x_i   \theta \sim \text{Bernoulli}(\theta)$	$\theta \sim \text{Beta}(\alpha, \beta)$	$\theta   x_{1:n} \sim \text{Beta}(\alpha + \sum_i x_i, \beta + \sum_i (1 - x_i))$
$x_i   \mu \sim \mathcal{N}(\mu, \sigma^2)$	$\mu \sim \mathcal{N}(\mu_0, 1)$	$\mu   x_{1:n} \sim \mathcal{N}\left(\frac{\sigma^2}{\sigma^2+n}(\mu_0 + \frac{1}{\sigma^2} \sum_i x_i), \frac{\sigma^2}{\sigma^2+n}\right)$
$x_i   \lambda \sim \text{Exponential}(\lambda)$	$\lambda \sim \text{Gamma}(\alpha, \beta)$	$\lambda   x_{1:n} \sim \text{Gamma}(\alpha + n, \beta + \sum_i x_i)$

### Generalized Linear Models

Regression	Inverse link function	Likelihood
Linear	identity	Gaussian
Logistic	sigmoid	Bernoulli
Poisson	exponential	Poisson
Negative binomial	exponential	Negative binomial

Sigmoid function:  $\sigma(x) = \frac{1}{1 + e^{-x}}$