

## Overview

Submit your writeup, including any code, as a PDF via gradescope.<sup>1</sup> We recommend reading through the entire homework beforehand and carefully using functions for testing procedures, plotting, and running experiments. Taking the time to reuse code will help in the long run!

Data science is a collaborative activity. While you may talk with others about the homework, please write up your solutions individually. If you discuss the homework with your peers, please include their names on your submission. Please make sure any handwritten answers are legible, as we may deduct points otherwise.

## Math Stats

1. (10 points) Work through the following exercises, and explain your reasoning in your answer.
  - (a) Suppose a particular drug test is 99% sensitive and 98% specific ([Here](#) is a Wikipedia link for a refresher on the terminology). The null hypothesis  $H_0$  is that the subject is not using the drug. Assume a prevalence of  $\pi_1 = 0.5\%$ , i.e. only 0.5% of people use the drug. Consider a randomly selected individual undergoing testing. Rounding to the nearest three significant figures, find
    - i. (1 point) the probability of testing positive given  $H_0$ .
    - ii. (1 point) the probability that they are not using the drug given they test positive.
    - iii. (2 points) the probability of testing positive a second time given they test positive once. You may assume the two tests are statistically independent given drug user status.
  - (b) Suppose we have a waiting time  $T \sim \text{Exponential}(\lambda)$  and wish to test

$$H_0 : \lambda = c \quad \text{vs} \quad H_1 : \lambda = 2c$$

for some  $c > 0$ . In this question, you'll use the *likelihood ratio test* (LRT) to compare these two hypotheses. The LRT considers the ratio of the two density functions  $f_1$  and  $f_0$  under the alternative and null respectively:

$$\text{LR}(T) = \frac{f_1(T)}{f_0(T)},$$

and rejects  $H_0$  when  $\text{LR}(T)$  is greater than some threshold  $\eta$ .

We use this test because of the *Neyman-Pearson lemma*, which states that the likelihood ratio test is the most powerful test (in other words, it has the highest power, or TPR) of significance level  $\alpha$ . That is, out of all possible tests of  $H_0$  vs  $H_1$  with  $\text{FPR} = \alpha$ , the likelihood ratio test has the highest TPR.

*Hint: For this question, you may find it helpful to brush up on computing probabilities involving continuous random variables. [Prob 140 textbook, Chapter 15](#) provides a helpful refresher.*

---

<sup>1</sup>In Jupyter, you can download as PDF or print to save as PDF

- i. (1 point) Compute  $LR(T)$  explicitly in terms of  $c$ .
- ii. (3 points) Let  $\alpha$  be our false positive rate ( $0 < \alpha < 1$ ). Compute the value of the threshold  $\eta$  so that the FPR of the test is equal to  $\alpha$ . We say that such a test has *significance level*  $\alpha$ . Your answer should be expressed in terms of  $\alpha$  and  $c$ .  
*Hint: start by expressing the FPR as a conditional probability, then connect it to the LRT decision rule and the densities  $f_0$  and  $f_1$ .*
- iii. (2 points) What is the *TPR* of this test? This is also known as the test's *power*. Your answer should be expressed in terms of  $\alpha$  and  $c$ .

## Bias in Police Stops

2. The following example is taken from [1, Ch. 6]:

A study of possible racial bias in police pedestrian stops was conducted in New York City in 2006. Each of  $N = 2749$  officers was assigned a score  $z_i$  on the basis of their stop data, with large positive values of  $z_i$  being possible evidence of bias. In computing  $z_i$ , an ingenious two-stage logistic regression analysis was used to compensate for differences in the time, place, and context of the individual stops.

We provide the data in a file `policez.csv`. Assume that under the null hypothesis (that the officers do not show racial bias), the  $z$ -scores should follow a standard normal distribution.

Note that throughout this question, the word “bias” refers to police officers’ racial bias, rather than the statistical term.

- (a) (1 point) In one plot, make a normalized histogram of the  $z$ -scores and a line plot of the pdf of the theoretical null  $\mathcal{N}(0, 1)$ . Describe how the fit looks.
- (b) (2 points) Compute  $p$ -values  $P_i = \Phi(-z_i)$  (where  $\Phi$  is the standard normal CDF) and then apply the BH procedure with  $\alpha = 0.2$ . Plot the sorted  $p$ -values as well as the decision boundary. How many discoveries did you make?
- (c) (2 points) Looking at the data, we can get a better fit to the distribution of  $z$ -scores if we use  $\mathcal{N}(0.10, 1.40^2)$ , called the empirical null (instead of the theoretical null from part (a)). Repeat steps (a) and (b), treating the empirical null as the null distribution.
- (d) (2 points) What assumption(s) are we implicitly making in part (c) by replacing the theoretical null  $\mathcal{N}(0, 1)$  with one which fits the data well  $\mathcal{N}(0.10, 1.40^2)$ ? What are the limitations of using the theoretical null? Which approach would you take when reporting discoveries of racial bias in this example? What other limitations do you see to this approach to modeling racial bias?

## $p$ -values, FDR and FWER

3. The `adult.csv` file contains data from a random sample of the US adult population. It includes two numerical fields: `Age` and `Hours worked per week`. It also includes four categorical fields (which we have binarized for you): `Gender`, `Education`, `Marriage status` and whether the person's income is greater than \$50,000. We will use this dataset to test the hypotheses of whether each of the categorical fields have any effect on the expectation of the numerical fields. For example, one test tests whether married individuals work significantly more or less than unmarried individuals.

- (a) (3 points) Write a function `avg_difference_in_means` that takes as input two column names: `binary_col`, the name of a column with binary data, and `numerical_col`, the name of a column with numerical data. The function should compute the  $p$ -value for a test of the following hypothesis test:

$H_0$  : There is no difference in the average value of `numerical_col` between the two groups specified in `binary_col`.

$H_1$  : The average value of `numerical_col` is different for the two groups specified in `binary_col`.

For example, the result of `avg_difference_in_means('Education', 'Age')` should be a  $p$ -value for testing whether there is a significant difference in age between college-educated and non-college-educated adults. You should use a permutation test (i.e., an A/B test from Data 8) to compute your  $p$ -values, using at least 25,000 permutations to form your final null distribution. Using such a large number of permutations will stabilize the  $p$ -values so that random noise is unlikely to lead to differing results across the class. On Datahub, running the full loop of tests should take a couple minutes.

*Hint:* It might be useful to recall how to run the simulations to get the necessary  $p$  values. [DATA 8 Textbook](#) provides a helpful refresher.

*Hint:* To shuffle a single column of a dataframe in pandas, you can use code similar to the following line. Make sure you use the correct arguments to the `sample` method!

```
df['my_column'] = df['my_column'].sample(...).values
```

- (b) (1 point) Use your function to compute eight  $p$ -values, one for each possible combination of categorical and numerical column.
- (c) (1 point) Suppose we use a naive  $p$ -value threshold of 0.05 to make a decision for each hypothesis test. Given the  $p$ -values from above, for which tests do we reject the null hypothesis?
- (d) (2 points) Suppose we want to guarantee a Family-wise Error Rate (FWER) of 0.05. Given the  $p$ -values from above, for which tests do we reject the null hypothesis?
- (e) (2 points) Suppose we want to guarantee a False Discovery Rate (FDR) of 0.05. Given the  $p$ -values from above, for which tests do we reject the null hypothesis?

*Hint:* Use the *Benjamini-Hochberg algorithm*.

- (f) (2 points) How do the results from (d) and (e) compare? Explain how and why these results are different.

*Hint: Recall how FWER and FDR are conceptually different.*

- (g) (2 points) Most variables don't always fit neatly into binary categories. As described earlier, we binarized these columns for you. Look at the original data in `adult_original.csv`. For one categorical column, give an example of how that variable could have been binarized differently, and how that might change the results from the earlier parts.

You aren't required to do any computation for this part: just explain how you might binarize one variable differently, and how that might change the results or your interpretation of them.

## References

- [1] Bradley Efron. *Large-scale inference: empirical Bayes methods for estimation, testing, and prediction*. Cambridge University Press, 2012.