

1. Let Y_{ij} be mutually independent random variables, $i = 1, \dots, A$ and $j = 1, \dots, n$, where $Y_{ij} \sim \text{beta}(\theta_i, 1)$. Recall on Day 1 you found that the asymptotic size α LRT of the hypothesis that $\theta_i = i\tau$ was

$$\phi(\mathbf{Y}) = \begin{cases} 1 & \text{if } TS(\mathbf{Y}) > \chi_{A-1, 1-\alpha}^2 \\ 0 & \text{otherwise} \end{cases}$$

where

$$TS(\mathbf{Y}) = -2n \left[\sum_{i=1}^A \ln(-\bar{\ell}_i) + \sum_{i=1}^A \ln(i\hat{\tau}) \right]$$

$$\bar{\ell}_i = \frac{\sum_{j=1}^n \ln(Y_{ij})}{n}$$

$$\hat{\tau} = \frac{-A}{\sum_{i=1}^A i\bar{\ell}_i}.$$

- (a) Using $A = 4$, $n = 3$, $\tau = 1$, and $\alpha = 0.05$ carry out a simulation study with 1000 replications to estimate the size of the test. Include your program and output.
- (b) What is the estimated size of the test? Does this test perform well? Explain how you arrived at your conclusion.
- (c) The simulation study used 1,000 replications. Is this a reasonable number of replications to estimate the size of a test? Explain why or why not.
2. A recent study shows that many of the pre-existing conditions that increase the mortality risk of COVID-19 are connected with long-term exposure to air pollution. In this problem, data on long term PM2.5 exposure ($\mu\text{g}/\text{m}^3$) and COVID-19 mortality counts are provided. The response is the cumulative COVID-19 mortality counts through May 12, 2020 for $n=3109$ US counties. County-level exposure to PM2.5 are calculated by averaging results from an established exposure prediction model for the years 2000–2016. In addition to PM2.5 exposure, 20 potential confounding variables are provided for the analysis.

The data files covid.Rdata (R dataset), covid.csv (CSV text), and covid.txt (tab delimited) all contain the following variables.

mean_pm25: county-level exposure to PM2.5
Deaths: cumulative COVID-19 mortality counts
Population: population

And some other variables that can be helpful:

| | |
|----------------------------------|---|
| <code>date_since:</code> | days since first COVID-19 case reported (a proxy for epidemic stage) |
| <code>q_popdensity:</code> | population density, |
| <code>older_percent:</code> | percent of population 65 years of age, |
| <code>prime_percent:</code> | percent of the population 45-64 years of age, |
| <code>mid_percent:</code> | percent of the population 15-44 years of age, |
| <code>poverty:</code> | percent of the population living in poverty, |
| <code>medhouseholdincome:</code> | median household income, |
| <code>pct_blk:</code> | percent of Black residents, |
| <code>hispanic:</code> | percent of Hispanic residents, |
| <code>education:</code> | percent of the adult population with less than high school education, |
| <code>medianhousevalue:</code> | median house value, |
| <code>pct_owner_occ:</code> | percent of owner-occupied housing, |
| <code>obese:</code> | percent of the population with obesity |
| <code>smoke:</code> | percent of current smokers, |
| <code>beds:</code> | number of hospital beds per unit population |
| <code>mean_summer_temp:</code> | average daily temperature for summer (June to September) |
| <code>mean_winter_temp:</code> | average daily temperature for winter (December to February) |
| <code>mean_summer_rm:</code> | relative humidity for summer (June to September) |
| <code>mean_winter_rm:</code> | relative humidity for winter (December to February) |
| <code>date_since_social:</code> | days since issuance of stay-at-home order for each state |

The researcher is interested in analyzing the connection between COVID deaths and long term PM2.5 exposure. You are a statistical consultant giving advice to the researcher.

You need to analyze the data and write a report that includes the following (note that your report must be written to a scientist who is only familiar with basic statistics):

- Description and the objectives of the study
- Fit two appropriate models for the count data to answer the researcher's question. For each analysis,
 - (a) Write out the model. Define each term and state model assumptions.
 - (b) Explain why the model is appropriate.
 - (c) Report the point and interval estimates of the parameter(s) that address the research question.
- The steps of the analysis and the findings. Make sure that you interpret the necessary results. Include plots and anything that is helpful to understand the analysis and results.
- Compare the analysis results from the two models you used. Explain which model is preferred and why.
- Draw a final conclusion to help the researcher understand the results of the statistical analysis.
- You also need to provide the SAR/R code in the Appendix.

Please note that the report has a strict three-page limit.