

Problem Set 1

Applied Stats II

Due: February 11, 2024

Instructions

- Please show your work! You may lose points by simply writing in the answer. If the problem requires you to execute commands in `R`, please include the code you used to get your answers. Please also include the `.R` file that contains your code. If you are not sure if work needs to be shown for a particular problem, please ask.
- Your homework should be submitted electronically on GitHub in `.pdf` form.
- This problem set is due before 23:59 on Sunday February 11, 2024. No late assignments will be accepted.

Question 1

The Kolmogorov-Smirnov test uses cumulative distribution statistics test the similarity of the empirical distribution of some observed data and a specified PDF, and serves as a goodness of fit test. The test statistic is created by:

$$D = \max_{i=1:n} \left\{ \frac{i}{n} - F_{(i)}, F_{(i)} - \frac{i-1}{n} \right\}$$

where F is the theoretical cumulative distribution of the distribution being tested and $F_{(i)}$ is the i th ordered value. Intuitively, the statistic takes the largest absolute difference between the two distribution functions across all x values. Large values indicate dissimilarity and the rejection of the hypothesis that the empirical distribution matches the queried theoretical distribution. The p-value is calculated from the Kolmogorov- Smirnov CDF:

$$p(D \leq d) = \frac{\sqrt{2\pi}}{d} \sum_{k=1}^{\infty} e^{-(2k-1)^2\pi^2/(8d^2)}$$

which generally requires approximation methods (see Marsaglia, Tsang, and Wang 2003). This so-called non-parametric test (this label comes from the fact that the distribution of the test statistic does not depend on the distribution of the data being tested) performs

poorly in small samples, but works well in a simulation environment. Write an R function that implements this test where the reference distribution is normal. Using R generate 1,000 Cauchy random variables (`rcauchy(1000, location = 0, scale = 1)`) and perform the test (remember, use the same seed, something like `set.seed(123)`, whenever you're generating your own data).

As a hint, you can create the empirical distribution and theoretical CDF using this code:

```

1  # create empirical distribution of observed data
2  ECDF <- ecdf(data)
3  empiricalCDF <- ECDF(data)
4  # generate test statistic
5  D <- max(abs(empiricalCDF - pnorm(data)))

1 # Set H0 and H1
2 # H0: There is no significant difference in the distribution of the two samples
3 # H1: There is a significant difference in the distribution of the two samples
4
5 # Set seeds to ensure reproducibility
6 set.seed(123)
7
8 # Generate 1000 Cauchy random variables
9 cauchy_data <- rcauchy(1000, location = 0, scale = 1)
10
11 # Create empirical distribution of observed data
12 ECDF <- ecdf(cauchy_data)
13 empiricalCDF <- ECDF(cauchy_data)
14
15 # Create theoretical cumulative distribution of observed data
16 theoretical_CDF <- pnorm(cauchy_data)
17 theoretical_CDF
18
19 # Calculate test statistic
20 D <- max(abs(empiricalCDF - theoretical_CDF))
21
22 # Calculate P value
23 p_value <- sqrt(2 * pi) / D * sum(exp(-(2*(1:1000)-1)^2 * pi^2 / (8 * D^2)))
24
25 # Print p value
26 print(p_value)
27
28 # p value is 5.652523e-29
29
30 # Check the results using ks.test()
31 results <- ks.test(empiricalCDF, theoretical_CDF)
32
33 # Print results
34 print(results)
35
36 # the results: D = 0.135, p-value = 2.432e-08, alternative hypothesis: two-
    sided

```

```
37 # so, there is a sufficient reason to reject H0,  
38 # and there is a significant difference in the distribution of the two samples
```

Question 2

Estimate an OLS regression in R that uses the Newton-Raphson algorithm (specifically BFGS, which is a quasi-Newton method), and show that you get the equivalent results to using `lm`. Use the code below to create your data.

```
1 # Set seeds to ensure reproducibility  
2 set.seed(123)  
3  
4 # Create dataframe  
5 data <- data.frame(x = runif(200, 1, 10))  
6 data$y <- 0 + 2.75*data$x + rnorm(200, 0, 1.5)  
7  
8 # Perform OLS linear regression using the lm() function  
9 lm_model <- lm(y ~ x, data = data)  
10  
11 summary(lm_model)  
12 # The intercept is 0.139187, the coefficient of x is 2.726699.  
13  
14 # Use the optim() function to perform linear regression  
15 # with the BFGS algorithm  
16 optim_model <- optim(par=c(0, 1),  
17                      fn=function(coef) {  
18                          sum((data$y - (coef[1] + coef[2] * data$x))^2)  
19                      }, method = "BFGS")  
20  
21 print(optim_model$par)  
22 # The intercept is 0.139187, the coefficient of x is 2.726699,  
23 # which is equivalent to using lm.
```