

Problem Set 1

Yuanyuan Liu

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:

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

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 # Function to perform the KS test
2 ks_test <- function(data) {
3   n <- length(data)
4   # Sort data for ECDF calculation
5   data <- sort(data)
6   ECDF <- ecdf(data)
7   empiricalCDF <- ECDF(data)
8
9   # Theoretical CDF for the standard normal distribution
10  theoreticalCDF <- pnorm(data)
11
12
13 # Function to calculate the p-value using matrix H
14 ks_pvalue_matrix <- function(D, n) {
15   k <- ceiling(n * D)
16   h <- k - n * D
17   m <- 2 * k - 1
18
19   # Fill the matrix
20   # Initialize the matrix H
21   H <- matrix(0, nrow = m, ncol = m)
22
23   for (i in 1:m) {
24     for (j in 2:i) {
25       H[i, j] <- (1-h^(i+1-j)) / factorial(i+1-j)
26     }
27   }
28
29   for (i in 1:(m-1)) {
30     for (j in 1) {
31       H[i, j] <- (1-h^i) / factorial(i)
32     }
33   }
34
35   for (i in 1:m) {
36     for (j in 1) {
37       H[i, j] <- (1-2*h^i) / factorial(i)

```

```

38     }
39   }
40
41   for (i in 2:(m-1)) {
42     for (j in 1:i) {
43       if (i+1 <= j) {
44         H[i, j] <- 1 / factorial(i+1-j)
45       } else {
46         H[i, j] <- 0
47       }
48     }
49   }
50
51   T <- H
52   for (i in 2:n) {
53     T <- T %*% H
54   }
55   p_value <- factorial(n) * T[k, k] / n^n
56   return(p_value)
57 }
58 D <- max(abs(empiricalCDF - theoreticalCDF))
59 p_value <- ks_pvalue_matrix(D, n)
60
61 return(list(D = D, p_value = p_value))
62 }
63
64 # Set seed for reproducibility
65 set.seed(123)
66
67 # Generate 1,000 Cauchy random variables
68 cauchy_data <- rcauchy(1000, location = 0, scale = 1)
69
70 # Perform the KS test
71 ks_result <- ks_test(cauchy_data)
72
73 # Print the result
74 print(ks_result)

```

The result is:
D:0.1347281
p_value:NaN

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
2 # Function to calculate the p-value using matrix H
3 ks_pvalue_matrix <- function(D, n) {

1 # Set the seed for reproducibility
2 set.seed(123)
3
4 # Create the data
5 data <- data.frame(x = runif(200, 1, 10))
6 data$y <- 0 + 2.75 * data$x + rnorm(200, 0, 1.5)
7
8 # Estimate the OLS regression using lm()
9 lm_fit <- lm(y ~ x, data = data)
10
11 # Now using optim() to perform OLS manually using the BFGS method
12 # Define the objective function (sum of squared residuals)
13 ssr <- function(params, data) {
14   with(data, sum((y - (params[1] + params[2] * x))^2))
15 }
16
17 # Initial parameter guesses
18 initial_params <- c(0, 0)
19
20 # Run optim() with BFGS method
21 optim_fit <- optim(par = initial_params, fn = ssr, data = data, method = "BFGS")
22
23 # Show the coefficients from lm and optim
24 lm_coefficients <- coef(lm_fit)
25 optim_coefficients <- optim_fit$par
26 # Print the results
27 print(lm_coefficients)
28 print(optim_coefficients)
```

Estimate an OLS regression using BFGS:

(Intercept) x

0.1391778 2.7267000

Estimate an OLS regression using `lm`:

(Intercept) x

0.1391874 2.7266985