

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 # Function to calculate the p-value using the fix formula
2 ks_p_value <- function(D) {
3   # fist part of the formula
4   pi_sqrt_2 <- sqrt(2 / pi)
5
6   # accumulate the sum of k
7   sum_val <- 0
8   k_max <- 500
9   for (k in 1:k_max) {
10     sum_val <- sum_val + exp(-2 * (2 * k - 1)^2 * pi^2 / (8 * D^2))
11   }
12
13   # Calculate the p-value using the accumulated sum
14   p_value <- pi_sqrt_2 * sum_val / D
15
16   # Ensure the p-value is within [0, 1]
17   p_value <- min(max(p_value, 0), 1)
18
19   return(p_value)
20 }
21
22 # Function to perform the Kolmogorov–Smirnov test
23 ks_test <- function(data) {
24
25   # Create empirical distribution of observed data
26   ECDF <- ecdf(data)
27   empiricalCDF <- ECDF(data)
28
29   # Generate test statistic D
30   D <- max(abs(empiricalCDF - pnorm(data)))
31
32   # Calculate the p-value using the formula function in last step and D
```

```

33  p_value <- ks_p_value(D)
34
35  # Return the test statistic and p-value
36  return(list(D = D, p_value = p_value))
37 }
38
39 set.seed(123)
40 # Generate 1,000 cauchy random variables
41 data <- rcauchy(1000, location = 0, scale = 1)
42
43 # Perform the Kolmogorov–Smirnov test
44 ks_test_result <- ks_test(data)
45
46 # Print the results
47 print(ks_test_result)
48 # Test Statistic (D): 0.1347281  p-value: 5.466419e-59
49
50 # valid with ks function
51 ks_result <- ks.test(data, "pnorm")
52
53 # Print the results
54 print(ks_result)
55 # Test Statistic (D) = 0.13573, p-value = 2.22e-16

```

Interpret result:

Both the custom K-S test and the built-in K-S test in R showed the same statistical conclusion: the sample data, generated from a Cauchy distribution, do not follow a normal distribution. Because The p-value is extremely small (5.466419×10^{-59} , 2.22×10^{-16} , though the p-value from build in ks test is not as small as the p-value from the custom test)

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 # data set generated
2 set.seed(123)
3 data <- data.frame(x = runif(200, 1, 10))
4 data$y <- 0 + 2.75*data$x + rnorm(200, 0, 1.5)
5
6 # step 1: function to calculate sum of square residual (ssr)
7 #data$x - Predictor, data$y - Response
8
9 ssr <- function(params) {
10   predicted_values <- params[1] + params[2] * data$x
11   residuals <- data$y - predicted_values
12   sum(residuals^2)
13 }
14
15 # step2: use the optim() function with the BFGS method to
16 #find the parameter estimates that minimize the sum of squared residuals
17
18 estimate_params <- c(0, 0) # guesses for intercept and slope
19 optim_results <- optim(par = estimate_params, fn = ssr, method = "BFGS")
20 bfgs_params <- optim_results$par
21
22 #For comparison, fit an OLS model using the lm() function
23 #Estimate Parameters using lm()
24 lm_model <- lm(y ~ x, data = data)
25 lm_params <- coef(lm_model)
26
27 print(bfgs_params)
28 #BFGS estimate for the intercept: 0.1391778 slope: 2.7267000
29
30 print(lm_params)
31 #lm() estimate for the intercept: 0.1391874 slope: 2.7266985
```

Interpret result:

By comparing findings from the output of `lm` and `bfgs`, there is a close match between the parameters estimated by the traditional `lm()` function and through using the BFGS method. The estimate intercept and slope values are nearly the same, with only a tiny difference in the last few decimal places. In general, both BFGS optimization method and `lm()` are accurate for estimating the parameters of an OLS regression, both methods able to calculate the best possible line that represents the relationship between X and Y in this data.