

# 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)))
```

Codes and test performing as below:

```
1 #####
2 # Problem 1
3 #####
4
5 ks_test_normal <- function(data){
6
7   #create empirical distribution of input data
8   ECDF <- ecdf(data)
9   empiricalCDF <- ECDF(data)
10
11  #create reference distribution CDF
12  normalCDF <- pnorm(data)
13
14  #generate statistic: largest absolute difference value
15  D <- max(abs(empiricalCDF - normalCDF))
16
17  #generate p value
18  p_value <- (sqrt(2*pi)/D) * sum(exp((-2*seq(1, length(data))-1)^2*pi^2)/(8*
19    D^2)))
20
21  #return the test statistic and p-value
22  return(list(D = D, p_value = p_value))
23 }
24
25
26 #generate Cauchy random variables
27 set.seed(123)
28 data <- rcauchy(1000, location = 0, scale = 1)
29
30 #execute the test
31 ks_test_normal(data)
32
33 #check the test
34 ks.test(data, "pnorm", mean = 0, sd = 1)
```

Output listed below:

```
$D[1] 0.1347281
$p_value[1] 5.652523e-29
```

Conclusion: as we can see from the results, if we set significant level as 0.05, p value is much less than it, which means the difference between the empirical distribution and theoretical distribution is enough significant to reject the H0: the observed data is from normal distribution.

## 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.seed(123)
2 data <- data.frame(x = runif(200, 1, 10))
3 data$y <- 0 + 2.75*data$x + rnorm(200, 0, 1.5)
```

Codes follow as below:

```
1 #calculate the RSS for an OLS regression
2 RSS_ols <- function(beta, x, y) {
3   y_hat <- beta[1] + beta[2] * x
4   sum((y - y_hat)^2)
5 }
6
7 #estimate the parameters using BFGS method
8 bfgs_result <- optim(fn = RSS_ols, par = 0:1, x = data$x, y = data$y, method =
   "BFGS")
9
10 #extract estimated coefficients after optimization
11 print(bfgs_result$par)
12
13 #get the equivalent result using lm()
14 lm_result <- lm(y ~ x, data)
15 print(coef(lm_result))
```

Output listed below:

```
> print(bfgs_result$par)
[1] 0.139187 2.726699

> print(coef(lm_result))
(Intercept)      x
0.1391874    2.7266985
```