

# 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:

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

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)))
```

### Code for problem 1

```
1 # Kolmogorov-Smirnov test function
2 set.seed(123)
3 kol_smir_test <- function(sample_size = 1000, location = 0, scale = 1) {
4   data <- rcauchy(sample_size, location = location, scale = scale)
5   ECDF <- ecdf(data)
6   empiricalCDF <- ECDF(data)
7   D <- max(abs(empiricalCDF - pnorm(data)))
8   print(paste("Kolmogorov-Smirnov test statistic (D) : ", D))
9   return(D)
10 }
11 # Calling the function
12 kol_smir_result <- kol_smir_test()
13 print(kol_smir_result)
14 # Calculating the p-value
15 D <- 0.13472806160635
16 pvalue_3 <- 1 - pnorm(sqrt(n) * D)
17 pvalue_3 <- 1 - pnorm(sqrt(1000) * 0.13472806160635)
18 print(pvalue_3)
19 # doing the ks.test to compare results
20 ks_result3 <- ks.test(data, "pnorm")
21 print(ks_result3)
```

**Output:** The p-value from my function is `pvalue-3 = 1.019963e-05` and the p-value from the `ks.test` is `ks-result3 = 2.22e-16`. I know mine is still higher than the `ks.test` but is the closest to zero I got after changing the function a lot. Actual output from `ks.test`: Asymptotic one-sample Kolmogorov-Smirnov test data: data D = 0.13573, p-value = 2.22e-16 alternative hypothesis: two-sided

## 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)
```

### Code

```
1 # Generating the data
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 # defining the function
6 OLS_obj <- function(beta, x, y) {
7   y_esp <- beta[1] + beta[2]*x
8   error <- y - y_esp
9   return(sum(error^2))
10 }
11 # Beta values at first
12 beta_um <- c(0, 0)
13 # Optimising with BFGS
14 bfg_otimo <- optim(par = beta_um, fn = OLS_obj, x = data$x, y = data$y, method
15   = "BFGS")
16 # getting the estimated coefficients
17 bfg_estcoef <- bfg_otimo$par
18 # Comparing the coefficients from BFGS with the lm
19 lm_comp <- lm(y ~ x, data = data)
20 coef_lm <- coef(lm_comp)
21 #
22 print(bfg_estcoef)
23 print(coef_lm)
```

**Output: bfg-estcoef = 0.1391778 , 2.7267000 and coef-lm = (Intercept) = 0.1391874 , x = 2.7266985**