

# 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)))
```

**My answer:**

```
1 ks_test_normal <- function(observed_data, seed = 123) {
2   set.seed(seed)
3
4   # Generate Cauchy random variables
5   cauchy_data <- rcauchy(1000, location = 0, scale = 1)
6
7   ECDF <- ecdf(observed_data)
8   empirical_CDF <- ECDF(observed_data)
9
10  theoretical_CDF <- pnorm(observed_data)
11
12  # Calculate test statistic
13  D <- max(abs(empirical_CDF - theoretical_CDF))
14
15  # Calculate p-value using Kolmogorov-Smirnov CDF
16  p_value <- sqrt(2 * pi) / D * sum(exp(-(2 * (1:1000) - 1)^2 * pi^2 / (8 * D
17    ^2)))
18
19  # Return the test statistic and p-value
20  return(list(test_statistic = D, p_value = p_value))
21 }
22 # Example usage:
23 observed_data <- rcauchy(1000, location = 0, scale = 1)
24 result <- ks_test_normal(observed_data, seed = 123)
25 print(result)
```

```
## $test_statistic
## [1] 0.1398682
##
## $p_value
## [1] 7.338799e-27
```

*I use this function, `ks_test_normal` to take observed data as input, generates Cauchy random variables, computes the empirical distribution, and then calculates the test statistic and*

*p-value. In this question, p-value calculation involves an infinite sum, so it might require numerical approximation methods.*

## 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 the random seed for reproducibility
2 set.seed(123)
3
4 # Create a data frame with two columns, where 'x' is generated from a uniform
   distribution in the range [1, 10],
5 # and 'y' is generated based on a linear relationship with added noise
6 data <- data.frame(x = runif(200, 1, 10))
7 data$y <- 0 + 2.75 * data$x + rnorm(200, 0, 1.5)
8
9 # Define the objective function for ordinary least squares (OLS)
10 ols_obj_function <- function(beta, x, y) {
11   y_pred <- beta[1] + beta[2] * x
12   residuals <- y - y_pred
13   sum(residuals^2)
14 }
15
16 # Use the BFGS optimization algorithm to fit the model with initial parameters
   set to c(0, 1)
17 result_optim <- optim(par = c(0, 1), fn = ols_obj_function, x = data$x, y =
   data$y, method = "BFGS")
18
19 # Get the parameter estimates from the optimization results
20 coef_optim <- result_optim$par
21
22 # Fit a linear model using the lm function, y ~ x
23 model_lm <- lm(y ~ x, data = data)
24
25 # Get the parameter estimates from the lm function
26 coef_lm <- coef(model_lm)
27
28 # Print the coefficients obtained using the BFGS optimization algorithm
29 cat("Coefficients using BFGS optimization:\n")
30 print(coef_optim)
31
32 # Print the coefficients obtained using the lm function
33 cat("\nCoefficients using lm function:\n")
34 print(coef_lm)
```

## The result of compiling:

```
## Coefficients using BFGS optimization:
```

```
print(coef_optim)
```

```
## [1] 0.139187 2.726699
```

```
cat("\nCoefficients using lm function:\n")
```

```
##  
## Coefficients using lm function:
```

```
print(coef_lm)
```

```
## (Intercept)          x  
## 0.1391874    2.7266985
```

*I use this function, `ols_obj` function that computes the sum of squared residuals for given coefficients. The `optim` function is then used to minimize this objective function with the BFGS algorithm. The results are compared with the `lm` function.*