

# Problem Set 1

## Applied Stats II

Due: February 12, 2023

### 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 12, 2023. 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)))
```

## Question 1 Answer

Using R, the following function was created:

```
1 ksfunc <- function(dat) {
2   # create empirical distribution of observed data
3   ECDF <- ecdf(dat)
4   empiricalCDF <- ECDF(dat)
5
6   # generate statistic
7   D <- max(abs(empiricalCDF - pnorm(dat)))
8
9   # generate p-value from base R function
10  res <- ks.test(empiricalCDF, pnorm(dat))
11
12  # generate p-value
13  k <- 1
14  fun <- function(n, x){
15    exp((( -((2*n)-1)^2)*(pi^2))/((8)*(x^2)))
16  }
17  mul <- ((sqrt(2*pi))/(D))
18  ans <- 0
19  tol <- 1e-06
20
21  while(TRUE) {
22    next_term <- fun(k, D)
23    next_term <- next_term*mul
24    ans <- ans + next_term
25    if(abs(ans)<tol){
26      break
27    }
28  }
```

```

28     k <- k+1
29   }
30
31   ans <- ans*mul
32   results <- data.frame(stat = D, p.value = res$p.value, func = ans)
33   return(results)
34
35 }

```

Using the function mentioned above to generate 1,000 Cauchy random variables, this data was then used to test for normality using the "ksfunc()" function and produced the following output:

```

stat      p.value      func
1 0.1347281 2.432277e-08 1.051657e-27

```

The "stat" value indicates the D statistic obtained from the first formula, while the "p.value" number indicates the P-value from the test obtained from R's built in ks.test() function. With the p-value being less than .05, we can reject the null hypothesis that the two CDF's are the same, meaning the data is not normal.

The "func" stat is meant to be the same as the "p.value" stat. Unfortunately, I could not implement the KS PDF function provided above into R to calculate the p-value. I tried to create the infinite sequence with a tolerance of (1e-06), but the first term seemed to be lower than the tolerance, so the function would stop running after the first term of the sequence.

## 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)

```

## Question 2 Answer:

First, the linear likelihood function was created, which finds the log of the likelihood of the PDF of the linear function, being

$$\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{y_i - x_i\beta}{2\sigma^2}}$$

This was done using the `linear1` function as shown below:

```
1 linear1 <- function(theta, y, X) {  
2   n <- nrow(X)  
3   k <- ncol(X)  
4   sigma2 <- theta[k+1]^2  
5   beta <- theta[1:k]  
6   e <- y - X%*%beta  
7   logl <- -.5*n*log(2*pi) - .5*n*log(sigma2) - ((t(e) %*% e) / (2*sigma2))  
8   return(-logl)  
9 }
```

The optimisation function in R was then used to find the MLE of the parameters, by utilising a Hessian matrix and the Newton-Raphson algorithm. The data was created using the code in the above question and passed through the `optim` function.

```
1 linear.MLE <- optim(fn=linear1 ,  
2                     par=c(1,1,1) ,  
3                     hessian = TRUE,  
4                     y = data_ex$y ,  
5                     X = cbind(1, data_ex$x) ,  
6                     method = "BFGS" )  
7 linear.MLE$par
```

This gave the following output:

```
[1] 0.1391893 0.2516985 1.4395471
```

An OLS linear regression was then run on the same data using `lm(data_ex$y ~ data_ex$x)`, which generated the following table: As shown, the X value, or  $\beta_0$  value, in the table corresponds to the second value in the output from the `linear.MLE` output, while the Constant, or  $\beta_1$  value corresponds to the first value of the output. This displays how OLS and MLS give the same output for these types of linear equations.

Table 1: Regression

	<i>Dependent variable:</i>
	y
x	0.252*** (0.042)
Constant	0.139 (0.253)
Observations	200
R <sup>2</sup>	0.156
Adjusted R <sup>2</sup>	0.152
Residual Std. Error	1.447 (df = 198)
F Statistic	36.629*** (df = 1; 198)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01