

# 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)))
```

My Answer for Question 1

```
1 set.seed(123) # as suggested in Question 1 to create reproducible results
2 # and so we get the same answers
3
4 # need to generate 1,000 Cauchy random variables
5 r <- 1000 # call this r
6 data1 <- rcauchy(r, location = 0, scale = 1)
7
8 # generate the Empirical Cumulative Distribution Function for the
9 # empirical distribution of observed data
10 ECDF <- ecdf(data)
11 empiricalCDF <- ECDF(data)
12
13 # generate the Test Statistic using the P value (p) we were given in R
14 D <- max(abs(empiricalCDF - pnorm(data)))
15 pi <- 3.142
16 p_value <- sqrt(2*pi)*sum((1)^2)*(pi)^2/(8*(D)^2)
17
18 # incorporated the given function
19 # Kolmogorov–Smirnov test (KS)
20
21 KS <- function(data) {
22   ECDF <- ecdf(data)
23   empiricalCDF < ECDF(data)
24   # create the Test Statistic
25   D <- max(abs(empiricalCDF - pnorm(data)))
26
27   sqrtroot <- sqrt(2*pi)
28   power_value <- ((1)^2)*((pi)^2)/(8*(D)^2)
29   p_value <- sqrtroot*sum(exp(power_value))
30   output <- list(D, p_value)
31   return(output)
32 }
33
34 ks.test(data1, "pnorm") # using the KS test R function
35
36
```

```

37 # RESULT
38 # One-sample Kolmogorov-Smirnov test
39 # data:  data1
40 # D = 0.13573, p-value = 2.22e-16
41 # alternative hypothesis: two-sided
42
43 # Outputs
44 # D = 0.13573
45 # p-value = 2.22e-16

```

## 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 # generate the Empirical Cumulative Distribution Function for the
2 # empirical distribution of observed data
3 ECDF <- ecdf(data)

```

My Answer for Question 2

```

1 # Given code
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 # I can estimate the Regression using the lm() function
7 coef(lm(data$y ~ data$x))
8 # this gives
9 # (Intercept)      data$x
10 # 0.1391874    2.7266985
11

```

```

12
13 # create log normal likelihood function
14 # using Tutorial 3 and the given code above to derive log-likelihood function
15
16 #
17 norm_likelihoood1 <- function(outcome, input, parameter) {
18   n <- nrow(input) # number of rows
19   k <- ncol(input)
20   beta <- parameter[1:k] # numbers or Betas or
21   sigma1 <- parameter[k+1]^2 # sigma helps to test the significance
22   e <- outcome - input%%beta
23   logl <- -.5*n*log(2*pi)-.5*n*log(sigma1) - ( (t(e) %% e) / (2*sigma1) )
24   return(-logl)
25 }
26
27 # generate the norm_likelihoood function and
28 # show that you get the equivalent results to using lm
29 # another way to generate the log likelihood function
30
31 norm_likelihoood2 <- function(outcome, input, parameter) {
32   n <- ncol(input)
33   beta <- parameter[1:n]
34   sigma <- sqrt(parameter[1+n])
35   -sum(dnorm(outcome, input %% beta, sigma, log=TRUE))
36 }
37
38 # print our estimated coefficients (intercept and beta_1)
39 # here the above functions norm_likelihoood1 and norm_likelihoood2
40 # can be run through the built in optim() function using the
41 # Newton Raphson algorithm which is an iterative procedure that can be used
42 # to calculate MLEs
43 results_norm1 <- optim(fn=norm_likelihoood1, outcome=data$y, input=cbind(1,
44   data$x), par=c(1,1,1), hessian=T, method="BFGS")
45
46 results_norm2 <- optim(fn=norm_likelihoood2, outcome=data$y, input=cbind(1,
47   data$x), par=c(1,1,1), hessian=T, method="BFGS")
48
49
50 results_norm1$par
51 # Intercept = 0.1396113
52 # Beta = 2.7266403
53
54 results_norm2$par
55 # Intercept = 0.1404966
56 # Beta = 2.7265080
57
58 # confirm that we get the same thing in with glm()
59 coef(lm(data$y~data$x))
60
61 # Results
62 # (Intercept)      data$x
63 # 0.1391874      2.7266985

```