

Problem Set 1

Applied Stats II

Due: February 14, 2022

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 class on Monday February 14, 2022. No late assignments will be accepted.
- Total available points for this homework is 80.

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).

Write an R function that implements this test where the reference distribution is normal.
Set the seed - `set.seed(123)`
Problem Set Broken down into 4 parts:

- 1) DATA GENERATION, 2) CALCULATE THEORETICAL AND EMPIRICAL CDFS,
- 3) CALCULATE D-STAT, 4) CALCULATE P-VALUE AND DRAW CONCLUSIONS

1) DATA GENERATION

Create 1000 Cauchy random variables and 1000 normal variables using `rcauchy` and `rnorm` functions

```
- x j- (rcauchy(1000, location = 0, scale = 1))  
- x  
- y j- rnorm(1000)  
- y  
Make dataframe from x and y  
- data j- data.frame(x, y)
```

2) FIND EMPIRICAL and THEORETICAL CDFS

Create empirical distribution of observed data using `ecdf()` function

```
- ECDF j- ecdf(data(dollar))x)  
- ECDF  
- empiricalCDF j- ECDF(data(dollar)x)  
- empiricalCDF
```

Create theoretical data with normal distribution using `ecdf()` function

- TCDF ;- ecdf(data(dollar)y) - TCDF - theoreticalCDF ;- TCDF(data(dollar)y) - theoreticalCDF

ALT method: Can also use pnorm(data, mean, sd) instead of ecdf for reference distribution

```
- mean(data(dollar)y) ANS = [1] 0.03609765
- sd(data(dollar)y) ANS = [1] 1.003827
- TCDF2 ;- pnorm(data(dollar)y, 0.03609765, 1.003827)
- TCDF2
- theoreticalCDF2 ;- TCDF(data(dollar)y)
- theoreticalCDF2
ANS = Same results as using ecdf() function
```

3) CALCULATE T-STAT NAMED D

```
- D ;- max(abs(empiricalCDF - pnorm(theoreticalCDF)))
- D ANS = [1] 0.8344237
```

4) CALCULATE P-VALUE + DRAW CONCLUSIONS

Install package 'dgof' to use ks.test function

```
- install.packages('dgof')
- library('dgof')
```

Try using ks.test function on theoreticalCDF with normal distribution

```
- ks.test(empiricalCDF, theoreticalCDF)
```

ANS = P-value = 1

Try using ks.test of empiricalCDFs normality using 'pnorm' as second argument

```
- ks.test(empiricalCDF, 'pnorm')
```

ANS = D = 0.5004, p-value ; 2.2e-16

P value far greater than 0.5, therefore indicating that ecdf is far from normal.

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.

First set the seed

```
- set.seed (123)
```

Create dataframe with x and y variables

```
- data <- data.frame(x = runif(200, 1, 10)) data$dollar.y <- 0 + 2.75*data$dollar.x +  
rnorm(200, 0, 1.5)
```

PART ONE: Try OLS method using `lm()` function

```
- lm <- lm(data$dollar.y ~ data$dollar.x)
```

```
- lm
```

ANS: [1] Coefficients: Intercept = -0.5055 Slope = 2.8166

PART TWO: Try alternate 'BFGS' method. ANS should be equivalent to `lm()`

Implement original function `minrss` to minimise the residual sum of squares

```
- minrss <- function(data, par) with(data, sum((par[1] + par[2] * x - y)^2))
```

Using `optim()` function following quasi-Newton BFGS method

```
- ?optim
```

```
- BFGSresult <- optim(par = c(0, 1), fn = min.RSS, data = data, gr=NULL, method =  
'BFGS', lower = -Inf, upper = Inf, control = list(), hessian = FALSE)
```

```
- BFGSresult
```

ANS = [1] -0.5055275 2.8165749

Coefficients using BFGS method equivalent to `lm()` function output.