

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)))
```

Answer:

An R function was written that implements the Kolmogorov-Smirnov test where the reference distribution is normal, using the below code:

```
1 # Create function to implement Kolmogorov-Smirnov test comparing data to
  normal distribution
2 ks_test <- function(data) {
3   # Create empirical distribution of observed data
4   ECDF <- ecdf(data)
5   empiricalCDF <- ECDF(data)
6   # Generate test statistic
7   D <- max(abs(empiricalCDF - pnorm(data)))
8   sum_val <- 0
9   for (i in 1:length(data)) {
10    sum_val <- sum_val + (exp(-(((2*i)-1)^2)*(pi^2)/((8*D)^2)))
11  }
12  p_val <- (sqrt(2*pi)/D) * sum_val
13  # Print results
14  print(cat("D =", D, "\n"))
15  print(cat("P-value =", p_val, "\n"))
16 }
```

1,000 Cauchy random variables were generated and the test was performed on this data, using the below code:

```
1 # Create data with Cauchy distribution
2 set.seed(2023)
3 data_emp <- rcauchy(1000, location=0, scale=1)

1 # Run function with our sample data
2 ks_test(data_emp)
```

This produced the below output:

D = 0.1320026
P-value = 0.002722055

This result was then checked against the built-in K-S test function in R, using the below code:

```
1 # Check results against built-in K-S test function
2 ks.test(data_emp, "pnorm")
```

This produced the below output:

Asymptotic one-sample Kolmogorov-Smirnov test

data: data_emp
D = 0.132, p-value = 1.443e-15
alternative hypothesis: two-sided

The two sets of results are functionally the same, with some small variation in the p-values probably due to rounding differences. The p-value indicates that we can reject the null hypothesis that the two samples were drawn from the same distribution.

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 # Create 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)
```

Answer:

After data was generated using the code above, the log of likelihood function was created as follows:

```
1 # Code log of likelihood function
2 linear.lik <- function(theta, y, X) {
3   n <- nrow(X)
4   k <- ncol(X)
5   beta <- theta[1:k]
6   sigma_sqrd <- theta[k+1]**2
7   e <- y - X%%beta
8   logl <- -0.5*n*log(2*pi) - 0.5*n*log(sigma_sqrd) -
9     ((t(e)%e)/(2*sigma_sqrd))
10  return(-logl)
11 }
```

The function was then ran using our data:

```
1 # Find parameters that maximize the function
2 linear.MLE <- optim(fn=linear.lik, par=c(1,1,1), hessian=TRUE,
3                    y=data$y, X=cbind(1, data$x), method = "BFGS")
4
5 linear.MLE$par
```

This gave the below output:

```
[1] 0.1398324 2.7265559 -1.4390716
```

We then ran the `lm()` function, producing the below output

```
1 # Find parameters that maximize the function
2 linear.MLE <- optim(fn=linear.lik, par=c(1,1,1), hessian=TRUE,
```

```

3 y=data$y, X=cbind(1, data$x), method = "BFGS")
4
5 linear.MLE$par

```

Table 1:

	<i>Dependent variable:</i>
	y
x	2.727*** (0.042)
Constant	0.139 (0.253)
Observations	200
R ²	0.956
Adjusted R ²	0.956
Residual Std. Error	1.447 (df = 198)
F Statistic	4,298.687*** (df = 1; 198)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

Comparing the results of our log of likelihood function to the output of the `lm()` function, we can see that the coefficient estimates are the same.