

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:

$$p(D \leq d) = \frac{\sqrt{2\pi}}{d} \sum_{k=1}^{\infty} e^{-(2k-1)^2\pi^2/(8d^2)}$$

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)))
```

Kolmogorov-Smirnov test function

```
1 # Create Kolmogorov-Smirnov test function
2 kolsmir <- function(data, p_value = 0.05){
3   # Sort input data
4   data_sorted <- sort(data)
5
6   # Calculate empirical CDF values
7   ecdf_values <- (1:length(data_sorted)) / length(data_sorted)
8
9   # Calculate theoretical CDF values
10  tcdf_values <- pcauchy(data_sorted, location = 0, scale = 1)
11
12  # Calculate maximum absolute difference between empirical and theoretical
13  kolsmir_statistic <- max(abs(ecdf_values - tcdf_values))
14
15  # Calculate p-value
16  kolsmir_p_value <- (sqrt(2 * pi)/kolsmir_statistic) * sum(exp(-((2 * nrow(
17    data_sorted)) - 1)^2 * pi^2 / (8 * kolsmir_statistic^2))))
18
19  # Confirm
20  kolsmir_result <- kolsmir_p_value < p_value
21
22  return(list("KS_Test_Statistic" = kolsmir_statistic, "KS_P-Value" = kolsmir_p_value, "KS_Test_Result" = kolsmir_result))
23 }
```

Implement custom function with the generated data. Test confirms TRUE.

```
1 # Implement Kolmogorov-Smirnov test function with generated data
2 data_result <- kolsmir(data = data)
3 print(data_result)
```

Lastly, Compare the given Reference values with the Custom Outputs

```
1 # Compare Reference and Custom Function Outputs
2 KS_compare <- data.frame("Custom Function" = data_result$KS_Test_Statistic, "
3   Reference" = D, "Difference" = abs(data_result$KS_Test_Statistic - D))
4 # Show how the Test Statistics compare between my custom function and the
5 print(KS_compare)
```

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 # Add y variable
2 data$y <- 0 + 2.75*data$x + rnorm(200, 0, 1.5)
3
4 # Use optim() function to estimate OLS regression
5 # Set starting values for the parameters as c(1, 1) to be optimised over
6 bfgs_model <- optim(par = c(1, 1),
```

The parameters are set to $c(1, 1)$ as a starting point for the algorithm to optimise over in order to reach the MLE. Note the *RSS* function is defined as a lambda function, though alternatively could be defined beforehand within the main or same function environment.

```
1
2 # Use optim() function to estimate OLS regression
3 # Set starting values for the parameters as c(1, 1) to be optimised over
4 bfgs_model <- optim(par = c(1, 1),
5                     # Set BFGS algorithm for the Newton–Raphson optimisation
6                     method = "BFGS",
7                     # Lambda function calculates the residual sum of squares
8                     # (ie differences between predicted and observed values)
9                     fn = function(coef) {
```

Then the built-in `lm()` function is used for comparison:

```
1
2 # Generate OLS regression using lm() for comparison
```

Now the coefficients are extracted for comparison:

```
1
2 # Extract coefficients for comparison
3 lm_coefficients <- lm_compare$coefficients
```

Lastly, demonstrate the coefficients are the same, confirming my custom model functions correctly.

```
1
2 # Compare outputs using all.equal() function in base R, ignoring differences
  in attributes such as column names.
3 model_check <- all.equal(lm_coefficients, bfgs_coefficients, check.attributes
  = FALSE)
```