

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)))
```

Answer 1

Rcauchy generates random deviates from the Cauchy. The length of the result is determined by n for rcauchy.

Reference: <https://stat.ethz.ch/R-manual/R-devel/library/stats/html/Cauchy.html>

```
1 # 1000 Cauchy random variables
2 cauchy_data <- rcauchy(1000, location = 0, scale = 1)
3 cauchy_data
```

```
1
2 kolmogorov-smirnov-test <- function(data) {
3   # Empirical distribution of observed data
4   ECDF <- ecdf(data)
5   empiricalCDF <- ECDF(data)
6
7   # Test statistic
8   D <- max(abs(empiricalCDF - pnorm(data)))
```

```

9
10 # P value calculation
11 p_value <- 2 * sum(exp(-(2*(1:length(data))-1)^2 / (8*D^2)))
12
13 return(p_value)
14 }

```

```

1
2 p_value <- kolmogorov_smirnov_test(cauchy_data)
3 p_value

```

p value is 0.002043117. If the p-value is low, it can be inferred that the two groups were drawn from populations that have dissimilar distributions.

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)

```

Answer 2

```

1 # OLS regression - the objective function
2 ols_obj_func <- function(beta, x, y) {
3   residuals <- y - beta[1] - beta[2] * x
4   sum(residuals^2)
5 }
6
7 # Estimate OLS regression with Newton-Raphson algorithm
8 result <- optim(c(0, 0), ols_obj_func, x = data$x, y = data$y, method = "BFGS")
9
10 # Estimated coefficients
11 estimated_coefs <- result$par
12
13 # Compare with lm function
14 lm_result <- lm(y ~ x, data = data)
15 lm_coefs <- coef(lm_result)
16
17 estimated_coefs
18 lm_coefs

```

```

estimated coefs
0.1391778 2.7267000
lm coefs

```

(Intercept) = 0.1391874

x = 2.7266985

The values are very close, with only slight differences in the decimal places. Overall, the estimated coefficients from the Newton Raphson algorithm and the `lm` function are more or less equal.