

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:

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

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)))
```

First, let's set up our null hypothesis and alternative hypothesis:

we want to know whether the empirical distribution (basically this is a 1,000 Cauchy random variables, generated by R) matches the queried theoretical distribution (in this question, this is a normal distribution). So, we can set up our H_0 and H_1 .

H_0 : The empirical distribution **follows** a normal distribution.

H_1 : The empirical distribution **does not follows** a normal distribution.

```
1 #####
2 # Problem 1
3 #####
4
5 # The function of Kolmogorov-Smirnov test with normal distribution
6 ks_normal <- function(data) {
7   # This function is used to test if a distribution follows normal distribution.
8
9   # Parameters:
10  # - data: The empirical distribution
11
12  # Returns:
13  # - D: largest absolute difference between the two distribution
14  # - p_value: Probability that the null hypothesis is correct
15
16  # create empirical distribution of observed data
17  ECDF <- ecdf(data)
18  empiricalCDF <- ECDF(data)
19  # generate test statistic
20  D <- max(abs(empiricalCDF - pnorm(data)))
21  # R can not deal with infinite, so I use 1000 instead
22  upper_limit <- 1000
23  # calculate p-value
24  p_value <- sqrt(2 * pi) / D * sum(exp(-(2 * (1:upper_limit) - 1)^2 * pi^2) / (8 *
    D^2)))
25  # return results
26  return(list(D = D, p_value = p_value))
27 }
28 # set seed as 123
29 set.seed(123)
30 # generate 1,000 Cauchy random variables
31 data <- rcauchy(1000, location = 0, scale = 1)
32 # use function to drive Kolmogorov-Smirnov test
33 ks_results <- ks_normal(data)
34 # print the test results
35 paste("D:", ks_results$D)
36 paste("p_value:", ks_results$p_value)
```

We got outputs from R:

```
[1] "D: 0.13472806160635"
[1] "p_value: 5.65252281681864e-29"
```

We can see that the p-value (= 5.65252281681864e-29) is below the $\alpha = 0.05$ threshold, so we would say that we find sufficient evidence to reject the null hypothesis that the empirical distribution follows a normal distribution.

So, our conclusion is to accept alternative hypothesis H_1 :

The empirical distribution **does not follows** a normal distribution.

We can also check answers by using the `ks.test` function in R:

```
1 # check by ks.test function in R
2 ks_check <- ks.test(data, "pnorm")
3 print(ks_check)
```

We got outputs from R:

Asymptotic one-sample Kolmogorov-Smirnov test

```
data: data
D = 0.13573, p-value = 2.22e-16
alternative hypothesis: two-sided
```

Compared with the function generated D (≈ 0.13573), our manually calculated D (≈ 0.13473) is nearly the same. And both p-value ($2.22e-16$ & $5.65252281681864e-29$) are very small and below the $\alpha = 0.05$ threshold.

Here is the explain of these 2 methods get the similar but not the same results:

- In `ks.test` command, R will use ∞ in formula. But in our handwritten formula, we replaced ∞ with a sufficiently large n, which is 1,000 in this case.

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 #####
2 # Problem 2
3 #####
4
5 set.seed(123)
6 data <- data.frame(x = runif(200, 1, 10))
7 data$y <- 0 + 2.75*data$x + rnorm(200, 0, 1.5)
```

Draw a scatter plot to see the relationship between the generated x and y.

```
1 # Draw a scatter to see the distribution
2 pdf("q2_plot1.pdf")
3 q2_scatter <- ggplot(data, aes(x = x, y = y)) +
4   geom_point(shape=1, size=2.8, color="#0F4C75") + # set point shape and color
5   theme_bw() + # set the coooooolest style
6   theme(panel.grid = element_blank()) # no grid
7 print(q2_scatter)
8 dev.off()
```

And we get this plot in R:

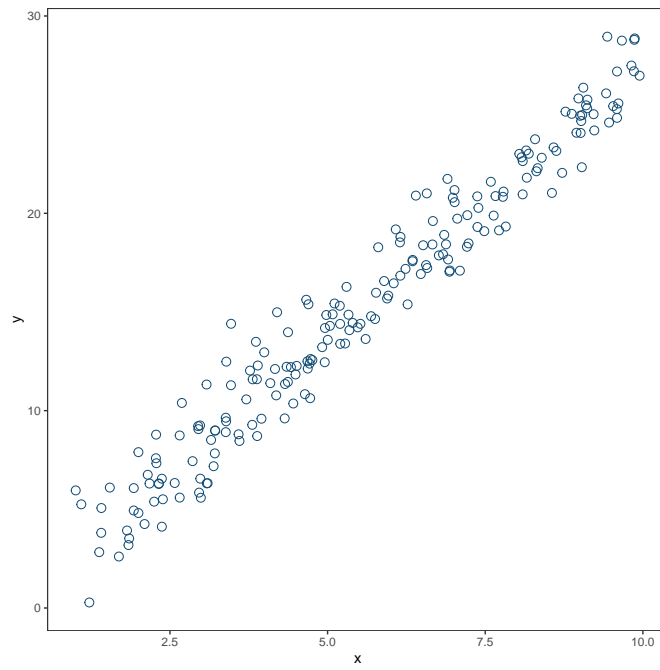


Figure 1: Scatter - Problem 2 data