

# Problem Set 1 Answers1

Kaley Burg

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)))
```

### My Answer:

- First I created my data

```
1 set.seed(123)

1 data <- (rcauchy(1000, location = 0, scale = 1))
```

- Then I created my Kolmogorov-Smirnoff CDF function:

```
1 kolsmir_pval <- function(data) {
2
3   set.seed(123)
4
5   # create empirical distribution of observed data
6
7   ECDF <- ecdf(data)
8   empiricalCDF <- ECDF(data)
9
10  #generating the test statistic
11
12  D <- max(abs(empiricalCDF - pnorm(data)))
13
14  #setting n
15  n <- length(data)
16
17  #initializing sum portion of the formula to 0
18  sum1 <- 0
19
20  #for loop for 1 to 1000 (instead of infinity)
21  for (k in 1:1000) {
22    sum1 <- sum1 + exp(-((2*k-1)^2)*(pi^2)/((8*D^2)))
23  }
24
25  #then we take the sum we created above and multiply by the other part
26  #of the formula
27}
```

```

28 #this gives the overall p value
29 p <- (((2*pi)^(1/2))/D)*sum1
30
31
32 #making it print like it would in ks.test
33
34 cat("D =", D, ", p-value =", p)
35
36 return(list(D = D, p = p))
37 }

```

- Lastly, I compared this with the built in ks.test function in R to check my answer

```
1 ks_result <- ks.test(data, pnorm)
```

- The results of both my function and the ks.test output are shown in the tables below:

Table 1: Kolmogorov-Smirnov Test - My function

Statistic	Value
$D$	0.1347281
$p$	$5.652523 \times 10^{-29}$

Table 2: Kolmogorov-Smirnov Test - ks.test Output

Statistic	Value
D	0.13573
p-value	$2.22 \times 10^{-16}$

- Note: the results are slightly different between my function and ks.test due to internal differences in the function, but they are similar enough to show that my function was implemented correctly.

#### **Interpretation:**

- As stated above, the D statistic is the largest absolute difference between the two distributions across all x values.
- As shown in the tables above, the D statistic is approximately 0.13 and the p-value associated with this statistic is very close to 0.
- This test overall assesses whether a sample comes from a specific distribution. Based on the p-statistic we obtain, we can draw the conclusion that we can **reject the null**

**hypothesis that the sample is drawn from the reference (normal) distribution.** This value would be significant at the  $\alpha < 0.001$  level, due to how small the obtained p-statistic is

## 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)
```

**My Answer:**

- First I ran the code given to us, which sets the seed and creates the data:

```
1 set.seed(123)
2 data2 <- data.frame(x = runif(200, 1, 10))
3 data2$y <- 0 + 2.75*data2$x + rnorm(200, 0, 1.5)
```

- Then I created a linear likelihood function to use in my Newton-Raphson algorithm (or more specifically, BFGS). This code came from the class slides and is utilizing the following function

$$\log(L) = \sum_{i=1}^n -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (y_i - x_i\beta)^2$$

```
1 linear_lik <- function(theta, y, X) {
2   n <- nrow(X)
3   k <- ncol(X)
4   beta <- theta[1:k]
5   sigma2 <- theta[k + 1]^2
6   e <- y - X%*%beta
7   logl <- -.5*n*log(2*pi) - .5*n*log(sigma2) - ((t(e) %*% e) / (2 * sigma2))
8   return(-logl)
9 }
```

- We assume that the points follow a normal probability distribution with mean  $x\beta$  and variance  $\sigma^2$ , i.e.,  $y \sim N(x\beta, \sigma^2)$
- So essentially, the code above is finding the parameters  $\beta$  and  $\sigma$  that maximize the probability for all points  $(x_i, y_i)$

- Then I used this function to find the maximum point on the surface created by the previous function using the optim function, using the BFGS algorithm to approximate the Newton-Raphson algorithm (this code came from the class slides):

```

1 linear_MLE <- optim(
2   fn = linear_lik,
3   par = c(1, 1, 1),
4   hessian = TRUE,
5   y = data2$y,
6   X = cbind(1, data2$x),
7   method = "BFGS"
8 )

```

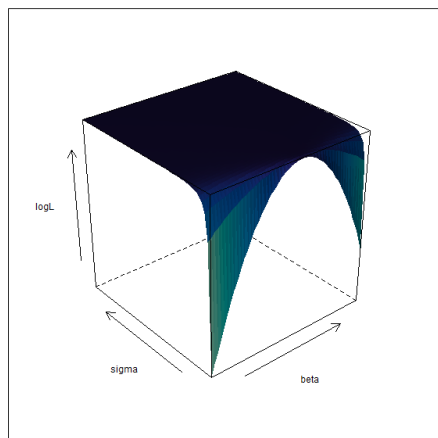
- For reference (not necessary for the overall answer/explanation) the surface itself can be visualized with this code (also from class slides):

```

1 surface <- list()
2 k <- 0
3 for (beta in seq(0, 5, 0.1)) {
4   for (sigma in seq(0.1, 5, 0.1)) {
5     k <- k + 1
6     logL <- linear_lik(theta = c(0, beta, sigma), y = data2$y, X = cbind
7       (1, data2$x))
8     surface[[k]] <- data.frame(beta = beta, sigma = sigma, logL = -logL)
9   }
10 }
11 surface <- do.call(rbind, surface)
12
13 library(lattice)
14
15 png(filename = "plot.png", width = 800, height = 600)
16 wireframe(logL ~ beta * sigma, surface, shade = TRUE)
17 dev.off()

```

- Which produces this plot:



- This function therefore produces an estimation of the parameters that specify this maximum point

Table 3: BFGS results

Parameter	Estimate
Intercept ( $\alpha$ )	0.1398324
$\beta$	2.7265559
$\sigma$	-1.4390716
$\sigma^2$	2.070927

- Then I compared the parameters created by my function with the results of the lm function

```

1 linear_MLE$par
2
3 lm_model <- summary(lm(y~x , data2))

```

Table 4: OLS (lm) results

Variable	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.13919	0.25276	0.551	0.582
x( $\beta$ )	2.72670	0.04159	65.564	< 2e - 16

- We can see in the highlighted terms that we are getting the same estimates for both the intercept ( $\alpha$ ) and  $\beta$  because both the BFGS algorithm and the OLS function are, in theory, consistent estimators of the true underlying parameter of interest

– In this case, an on average effect of the covariate on the outcome

- At some point the answers we get using the different methods would not be the same because they are different estimators
- For instance, OLS is closed form in the sense that we have the data and we plug it in to get estimates, while MLE is more complex as it uses numeric approximations
- So, in theory if we had a large amount of data, they will get closer to true estimate, but MLE is **not an unbiased estimator**