

Answer Key: Problem Set 1

Applied Stats II

Jeffrey Ziegler

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 class on Monday February 14, 2022. No late assignments will be accepted.*
- *Total available points for this homework is 80.*

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-Smirnoff 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, we'll create our data using `rcauchy(1000, location = 0, scale = 1)`. Then, we can insert the above code into our function, which first calculates our test statistic and then iteratively calculates our p-value.

```
1 # create data
2 set.seed(5)
3 n <- 1000
4 empirical <- rcauchy(n, location = 0, scale = 1)
5 # create K-S function to compare data against normal distribution
6 ksTest <- function(data){
7   # create empirical distribution of observed data
8   ECDF <- ecdf(data)
9   empiricalCDF <- ECDF(data)
10  # generate test statistic
11  D <- max(abs(empiricalCDF - pnorm(data)))
12  # calculate p-value
13  # empty vector to be filled
14  summed <- NULL
15  for(i in 1:n){
16    summed <- c(summed, exp((-2 * i - 1)^2 * pi^2) / ((8 * D)^2))
17  }
18  pValue <- sqrt(2*pi)/D * sum(summed)
19  cat("D =", D, "\n")
20  cat("p-value =", pValue, "\n")
21 }
22 # run "by hand" K-S test
23 ksTest(empirical)
```

```
1 D = 0.135986
2 p-value = 0.004404211
```

We can corroborate that the above output from our function operates similarly to the built-in function `ks.test()`, though there are some differences due to approximation and rounding.

```
1 # double check K-S test with R function
2 ks.test(empirical, "pnorm")
```

```

1 One-sample Kolmogorov-Smirnov test
2
3 data: empirical
4 D = 0.13699, p-value < 2.2e-16
5 alternative hypothesis: two-sided

```

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 ex_data <- data.frame(x = runif(200, 1, 10))
3 ex_data$y <- 0 + 2.75*ex_data$x + rnorm(200, 0, 1.5)

```

First, we need to create our data, then we create write our log-likelihood function.

```

1 # create log-likelihood function
2 # we need to specify what the outcome will be,
3 # what the input variables are, and the starting values
4 # of the parameters to be estimated
5 norm_log_likelihood <- function(outcome, input, parameter) {
6   # how many betas are we estimating? (# of columns)
7   n <- ncol(input)
8   beta <- parameter[1:n]
9   # estimate of our variance as well
10  sigma <- sqrt(parameter[1+n])
11  # put all of those into our log-likelihood (since log=T)
12  # remember, we're referencing normal distribution, so
13  # that's why you see dnorm()
14  -sum(dnorm(outcome, input %*% beta, sigma, log=TRUE))
15 }

```

Next, we can execute our function and confirm that we get the same answer using `lm()`.

```

1 # print our estimated coefficients (intercept and beta_1)
2 results_norm <- optim(fn=norm_log_likelihood, outcome=ex_data$y, input=cbind
3   (1, ex_data$x), par=c(1, 1, 1), hessian=T)
4 # print our estimated coefficients (intercept and beta_1)
5 # get same results regardless of which log-likelihood function we use
6 round(results_norm$par, 2)[1:2]

```

```

1 [1] 0.14 2.73

```

```

1 # confirm that we get the same thing with glm()
2 round(coef(lm(y~x, data=ex_data)), 2)

```

```

1 [1] 0.14 2.73

```