

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)))
```

Set up null hypothesis(H0) and alternative hypothesis(H1):

Null Hypothesis (H0): The observed data follows a specified theoretical distribution (in this problem, the specified distribution is the standard normal distribution).

Alternative Hypothesis (H1): The observed data does not follow the specified theoretical distribution.

```
1 # Define function for Kolmogorov-Smirnov test with normal reference
  distribution
2 kolmogorov_smirnov_test <- function(sample_size, seed = 123) {
3   # Set seed for reproducibility
4   set.seed(seed)
5
6   # Generate Cauchy random variables, where sample_size is 1000
7   data <- rcauchy(sample_size, location = 0, scale = 1)
8
9   # Create empirical distribution of observed data
10  ECDF <- ecdf(data) # Create empirical distribution function of observed
    data
11  empiricalCDF <- ECDF(data) # Compute empirical distribution function values
    of observed data
12
13  # Generate test statistic
14  D <- max(abs(empiricalCDF - pnorm(data))) # Generate test statistic, which
    is the maximum difference between the empirical distribution function of
    the observed data and the theoretical cumulative distribution function of
    the standard normal distribution
15
16  # Calculate p-value using Kolmogorov-Smirnov CDF
17  p_value <- sqrt(2 * pi) / D * sum(exp(-((2 * (1:length(data)) - 1)^2 * pi^2)
    / (8 * D^2)))
18
```

```

19 # Return test statistic and p-value
20 return(list(D = D, p_value = p_value))
21 }
22
23 # Example usage
24 result <- kolmogorov-smirnov-test(1000)
25 print(result)
26
27 # $D
28 # [1] 0.1347281
29
30 # $p_value
31 # [1] 5.652523e-29
32
33
34 # Compare with ks.test() function in R to validate custom implementation
35 ks_result <- ks.test(data, "pnorm")
36 print(ks_result)
37
38 # Asymptotic one-sample Kolmogorov-Smirnov test
39 # data: data
40 # D = 0.13573, p-value = 2.22e-16
41 # alternative hypothesis: two-sided

```

Conclusion: The p-value obtained from the Kolmogorov-Smirnov test is extremely small (below $= 0.05$, and close to zero), indicating strong evidence against the null hypothesis. Therefore, we reject the null hypothesis and conclude that the observed data does not follow the specified theoretical distribution (standard normal distribution).

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.

Solution Approach:

1. First, generate the dataset.
2. Define the negative log-likelihood function, which is used to minimize the sum of squared residuals during fitting.
3. Optimize the negative log-likelihood function using the BFGS algorithm to obtain the estimates of regression coefficients.
4. Use the `lm()` function to perform ordinary least squares linear regression and obtain the estimates of regression coefficients.
5. Compare the coefficients obtained from both methods to validate their equivalence.

```

1 # Step 1: Generate data

```

```

2 set.seed(123)
3 # Create dataset
4 data <- data.frame(x = runif(200, 1, 10))
5 data$y <- 0 + 2.75 * data$x + rnorm(200, 0, 1.5)
6
7 # Define negative log-likelihood function
8 linear.lik <- function(theta, y, X) {
9   n <- nrow(X) # Number of observations
10  k <- ncol(X) # Number of predictors
11  beta <- theta[1:k] # Regression coefficients
12  sigma2 <- theta[k+1]^2 # Variance
13  e <- y - X %*% beta # Residuals
14
15  # Negative log-likelihood function
16  logl <- -0.5 * n * log(2 * pi) - 0.5 * n * log(sigma2) - sum((t(e) %*% e) /
17    (2 * sigma2))
18  return(-logl)
19 }
20 # Optimize negative log-likelihood function using BFGS algorithm
21 opt_result <- optim(par = c(0, 0, 1), fn = linear.lik, y = data$y, X = cbind
22   (1, data$x), method = "BFGS")
23 # Extract estimated coefficients from BFGS algorithm
24 bfgs_coefs <- optim_result$par
25
26 # Print coefficients from BFGS algorithm
27 cat("Coefficients from BFGS algorithm:")
28 print(bfgs_coefs)
29 # 0.139187 2.726699
30
31 # Estimated coefficients using BFGS algorithm
32 beta_hat <- opt_result$par
33 cat("Estimated coefficients using BFGS algorithm:", beta_hat)
34 # Estimated coefficients using BFGS algorithm: -0.4065592 -0.4595849 438.0535
35
36 # Compare with lm() function, Ordinary Least Squares Linear Regression
37 summary(lm(y ~ x, data))
38
39 #
40 # Call:
41 # lm(formula = y ~ x, data = data)
42 #
43 # Residuals:
44 #   Min       1Q   Median       3Q      Max
45 # -3.1906 -0.9374 -0.1665  0.8931  4.8032
46
47 # Coefficients:
48 #   Estimate Std. Error t value Pr(>|t|)
49 # (Intercept)  0.13919    0.25276   0.551   0.582
50 # x            2.72670    0.04159  65.564 <2e-16 ***

```

```

51 #      ———
52 #   Signif. codes:  0      ***      0.001      **      0.01      *      0.05      .      0.1
53                      1
54 # Residual standard error: 1.447 on 198 degrees of freedom
55 # Multiple R-squared:  0.956, Adjusted R-squared:  0.9557
56 # F-statistic:  4299 on 1 and 198 DF,  p-value: < 2.2e-16

```

Conclusion: Based on the summary statistics of the regression model, the linear regression model performs well in explaining the dependent variable. The estimated coefficients are statistically significant as indicated by their very small p-values. The high value of multiple R-squared, close to 1, indicates that the model can explain a large proportion of the variance in the dependent variable. The p-value of the F-statistic, approaching zero, suggests that the overall model fit is significant.

The OLS regression coefficients estimated using the Newton-Raphson algorithm with the BFGS optimization method are very close to those obtained using the `lm()` function, indicating their equivalence.