

Jason_McCoy_Extra Credit Problem #4 (15 points)

Fisher proposed a method for calculating a confidence interval for the Pearson Correlation Coefficient. This method involves transformation of r , the sample Pearson Correlation Coefficient. A z-score is calculated and confidence interval is constructed. This confidence interval is then transformed back to the original scale of measurement. This method is explained in the links:

<http://www2.sas.com/proceedings/sugi31/170-31.pdf>

http://onlinestatbook.com/2/estimation/correlation_ci.html

Use the data provided and construct the data frame “test” with the code below. The data frame contains test results for 49 students on two standardized tests. Each student took both tests. Do not change the order of the two test entries or the matching per student will not be correct

```
testA <- c(58,49.7,51.4,51.8,57.5,52.4,47.8,45.7,51.7,46,50.4,61.9,49.6,61.6,54,54.9,49.7,
          47.9,59.8,52.3,48.4,49.1,53.7,48.4,47.6,50.8,58.2,59.8,42.7,47.8,51.4,50.9,49.4,
          64.1,51.7,48.7,48.3,46.1,47.3,57.7,41.8,51.5,46.9,42,50.5,46.3,44,59.3,52.8)
testB <- c(56.1,51.5,52.8,52.5,57.4,53.86,48.5,49.8,53.9,49.3,51.8,60,51.4,60.2,53.8,52,
          49,49.7,59.9,51.2,51.6,49.3,53.8,50.7,50.8,49.8,59,56.6,47.7,47.2,50.9,53.3,
          50.6,60.1,50.6,50,48.5,47.8,47.8,55.1,44.9,51.9,50.3,44.3,52,49,46.2,59,52)

test <- as.data.frame(cbind(testA,testB))

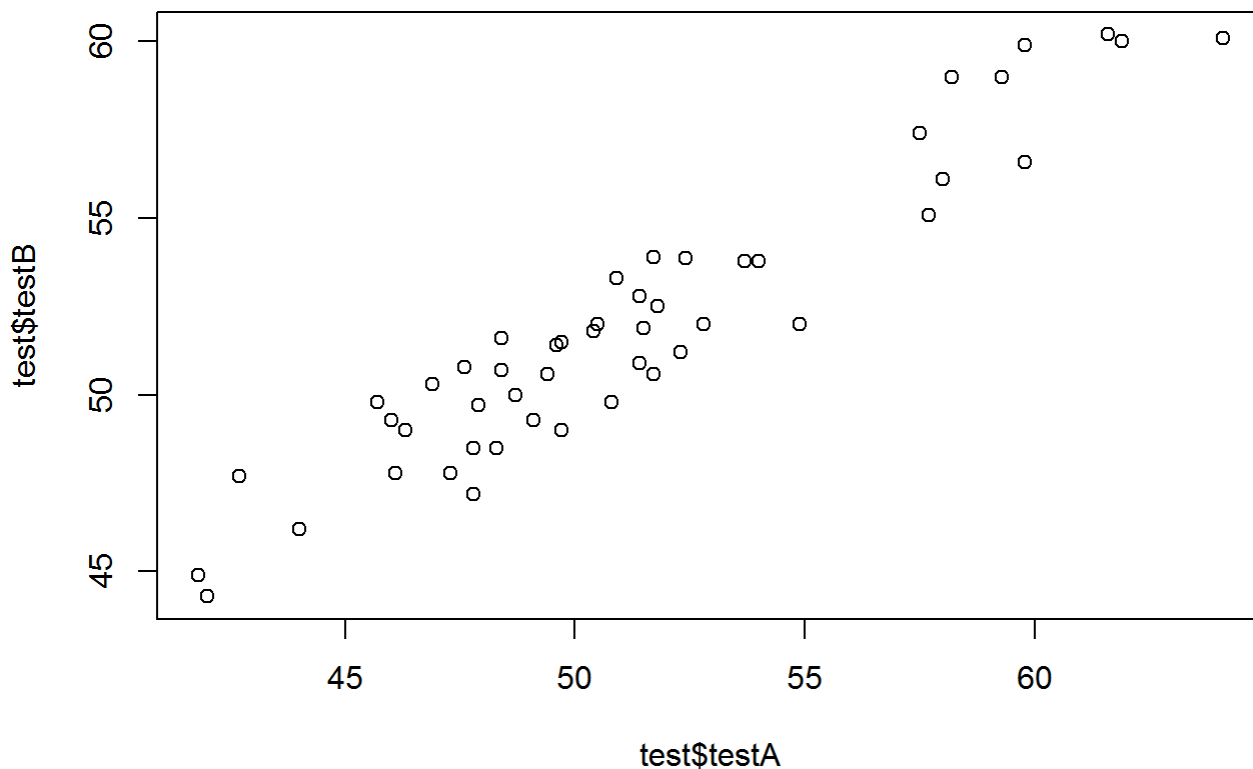
str(test)
```

```
## 'data.frame':    49 obs. of  2 variables:
##  $ testA: num  58 49.7 51.4 51.8 57.5 52.4 47.8 45.7 51.7 46 ...
##  $ testB: num  56.1 51.5 52.8 52.5 57.4 ...
```

```
summary(test)
```

```
##           testA           testB
##  Min.      :41.80   Min.      :44.30
##  1st Qu.:47.80   1st Qu.:49.30
##  Median :50.50   Median :51.40
##  Mean   :51.25   Mean   :51.95
##  3rd Qu.:53.70   3rd Qu.:53.80
##  Max.    :64.10   Max.    :60.20
```

```
plot(test$testA,test$testB)
```



Part #1 (3 points)

Determine a 95% confidence interval for the Pearson Correlation Coefficient of the data in “test” using Fisher’s method. Present the code and the confidence interval for rho, the Pearson Correlation Coefficient. Calculations can be simplified using *tanh()* and *atanh()*. Also, submit the data to *cor.test()* and present those results as well.

```
rho = cor(testA, testB)
zr = 0.5*log((1+rho)/(1-rho))
zcrit = qnorm(0.975, mean = 0, sd = 1, lower.tail = TRUE, log.p = FALSE)
cl = zr - zcrit*sqrt(1/(length(testA)-3))
cu = zr + zcrit*sqrt(1/(length(testA)-3))

#Confidence interval limits
rho_lower = (exp(2*cl)-1)/(exp(2*cl)+1)
rho_upper = (exp(2*cu)-1)/(exp(2*cu)+1)
rho_lower
```

```
## [1] 0.9113003
```

```
rho_upper
```

```
## [1] 0.9712052
```

```
#cor.test results
```

```
cor.test(test$testA, test$testB, alternative = "two.sided", method = "pearson", exact = NULL, conf.level = 0.95, continuity = FALSE)
```

```
##
## Pearson's product-moment correlation
##
## data: test$testA and test$testB
## t = 20.69, df = 47, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## 0.9113003 0.9712052
## sample estimates:
## cor
## 0.9492478
```

```
#The results match exactly!
```

Part #2 (5 points)

Bootstrapping can be used to arrive at an estimated confidence interval. The process involves resampling with replacement of rows from “test.” The first step is to randomly sample with replacement the 49 rows of “test”. Each sample will consist of 49 rows for which a sample correlation coefficient is calculated. For this purpose, the function *sample.int()* may be used to determine a sample of row numbers to be used. The function *cor()* should be used to determine the sample correlation coefficient. This step is repeated 10,000 times resulting in 10,000 sample correlation coefficients. The 10,000 calculated sample correlation coefficients are written to a vector. *quantile()* is passed this vector to calculate the 2.5% (0.025) and 97.5% (0.975) quantiles which determines the 95% Percentile Bootstrap confidence interval.

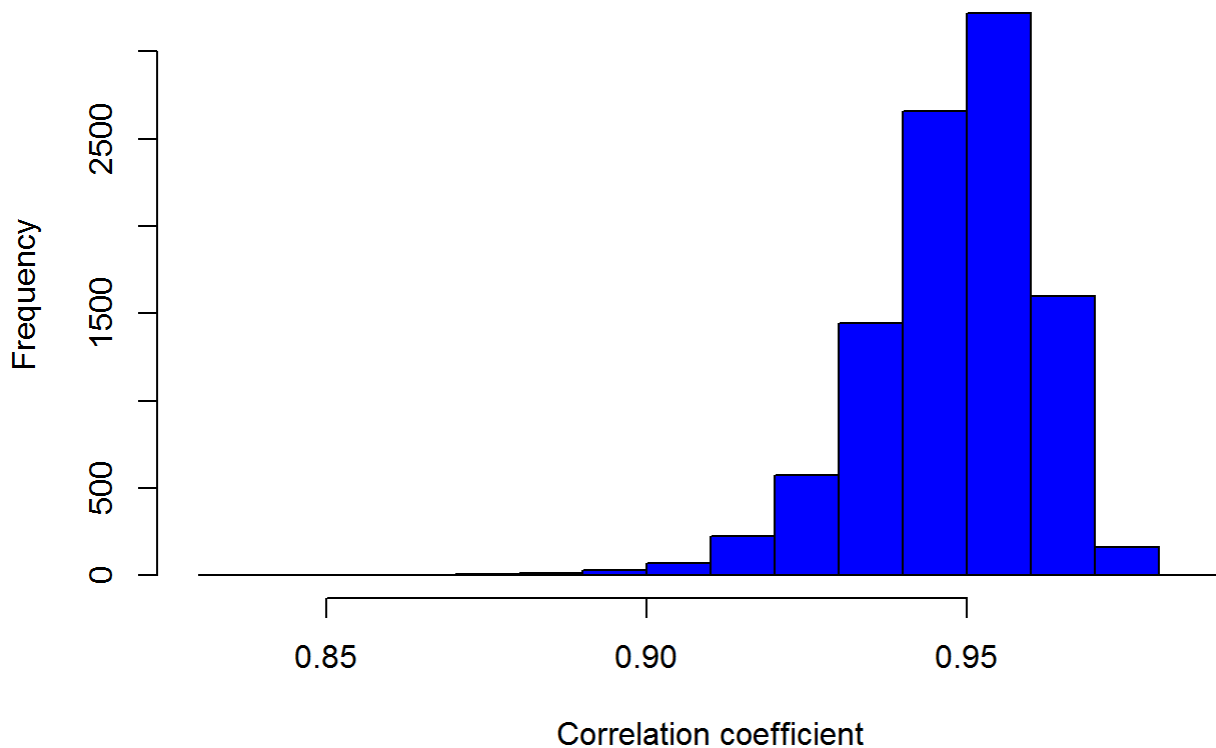
Refer to the course site library reserves and read the listing for Section 9.5 of “Mathematical Statistics with Resampling and R,” by Chihara and Hesterberg.

You will write code which does this work. Use *set.seed(123)*. A “for” loop may be used to repeat the sampling and correlation coefficient calculations. Plot a histogram and report a two-sided 95% confidence interval.

```
set.seed(123)
sample_corr = c()
for (i in 1:10000) {
  x=test[sample.int(nrow(test),49,replace=TRUE),]
  sample_corr[i] = cor(x$testA, x$testB)
}

#Histogram of correlations
hist(sample_corr, xlab="Correlation coefficient", main="Histogram of correlation coefficients", col="blue")
```

Histogram of correlation coefficients



```
#Confidence Interval
quantile(sample_corr, probs=c(0.025, 0.975))
```

```
##      2.5%      97.5%
## 0.9176197 0.9687603
```

How do bootstrapping results compare to the results from Part 1?

Answer: Bootstrapping results are almost similar to those in Part 1. The confidence interval obtained from Bootstrapping is slightly narrower than that obtained in Part 1.

Part #3 (7 points)

Bootstrapping can also be used to arrive at confidence intervals for regression coefficients. This process is similar to the process in Part #2. Using the current data frame, rows of “test” are randomly sampled with replacement. Each sample is passed to `*lm()` and the coefficients extracted. Section 9.5 of Chihara and Hesterberg give an example R script showing how this may be accomplished.

Write code using “test” to produce histograms and 95% two-sided bootstrap confidence intervals for the intercept and slope of a simple linear regression. Please keep `set.seed(123)`. A “for” loop can be written to sample via `sample.int()`, a linear regression model fit via `lm()` and the coefficients extracted via `coef()`. Note that we pass our fitted model object to `coef()` to return the coefficients; e.g. `coef(model)[1]` should return the intercept, `coef(model)[2]` the slope.

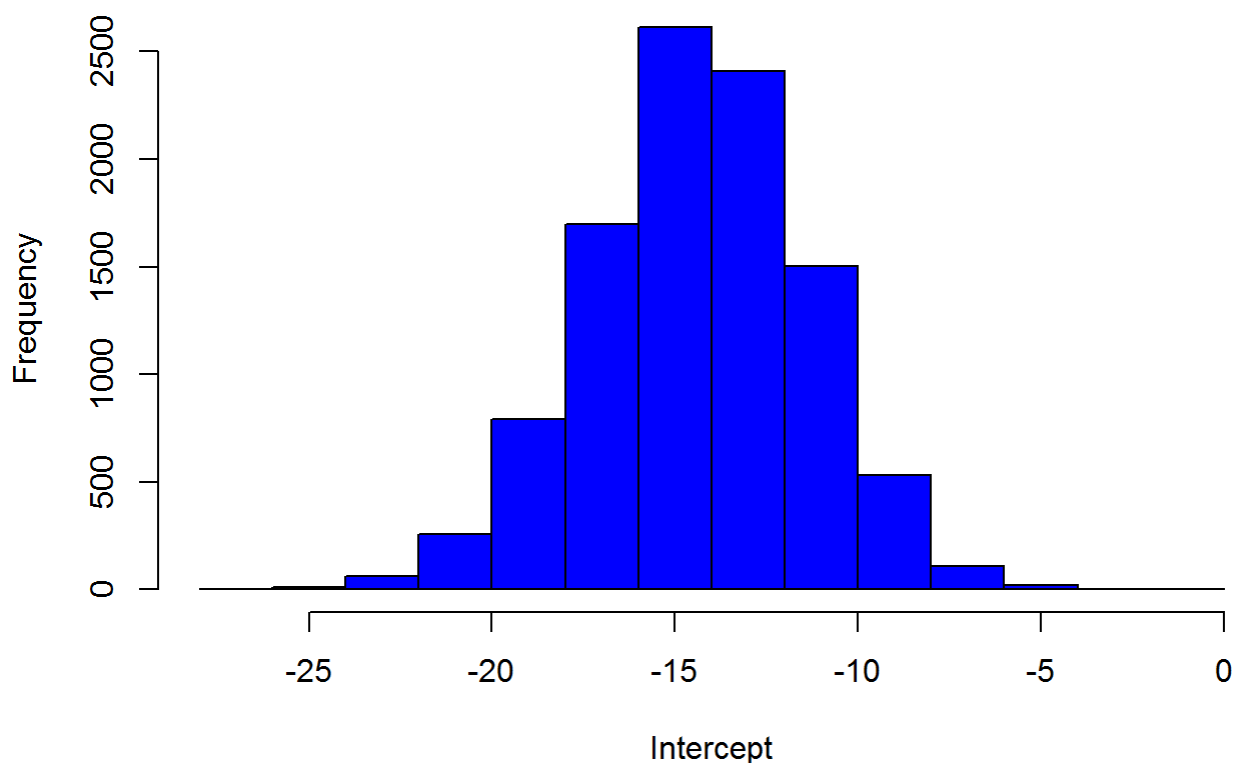
Present two histograms, one for bootstrap intercept results, and one for the bootstrap slope results showing the 2.5% and

97.5% quantiles on each. In addition, show the location on the histograms of the estimated intercept and slope of test using `lm()` without bootstrapping.

Lastly, generate a scatter plot of the estimated bootstrap slopes versus the estimated bootstrap intercepts. There will be 10,000 points appearing in this plot, one for each bootstrap sample. Place the intercept on the x-axis and the slope on the y-axis.

```
set.seed(123)
intercept =c()
slope =c()
for (i in 1:10000) {
  x=test[sample.int(nrow(test),49,replace=TRUE),]
  model=lm(x$testA~x$testB)
  intercept[i]=coef(model)[1]
  slope[i]=coef(model)[2]
}
hist(intercept, xlab="Intercept",main="Histogram of Intercepts",col="blue")
```

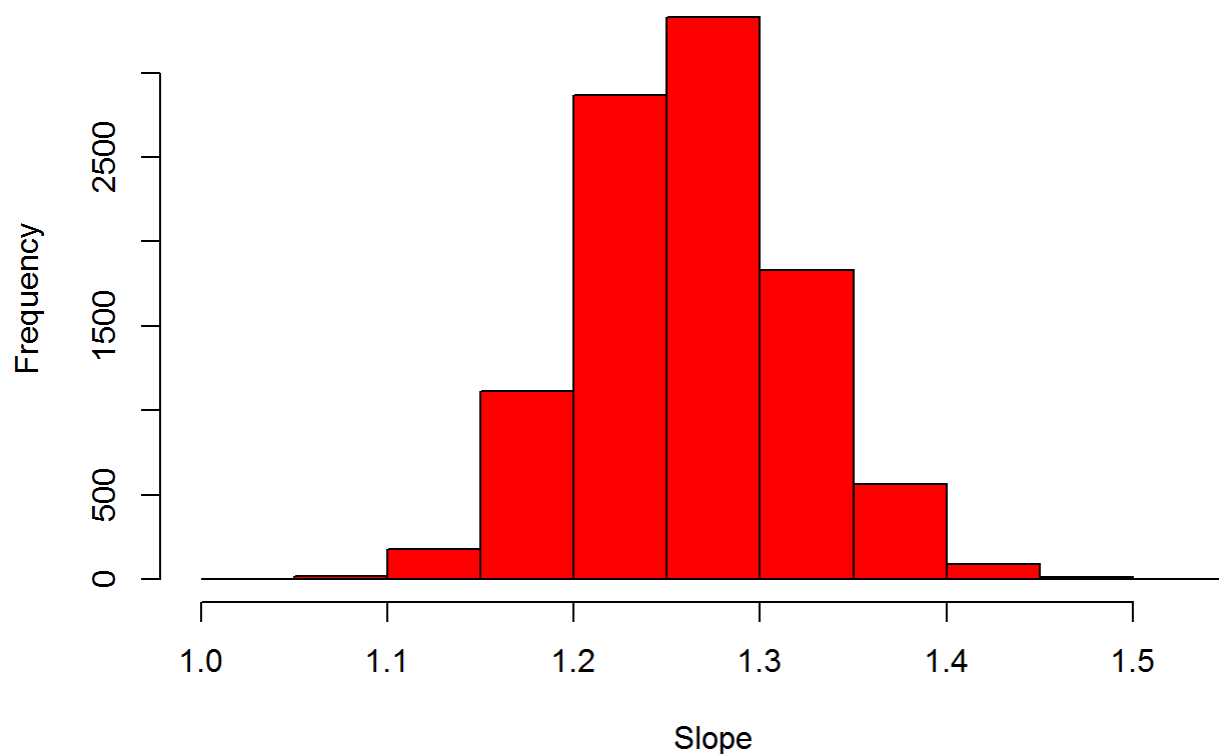
Histogram of Intercepts



```
hist(slope, xlab="Slope",main="Histogram of Slope",col="red")
```

-2 points Note the highlighted statement above. Estimated values are not shown.

Histogram of Slope



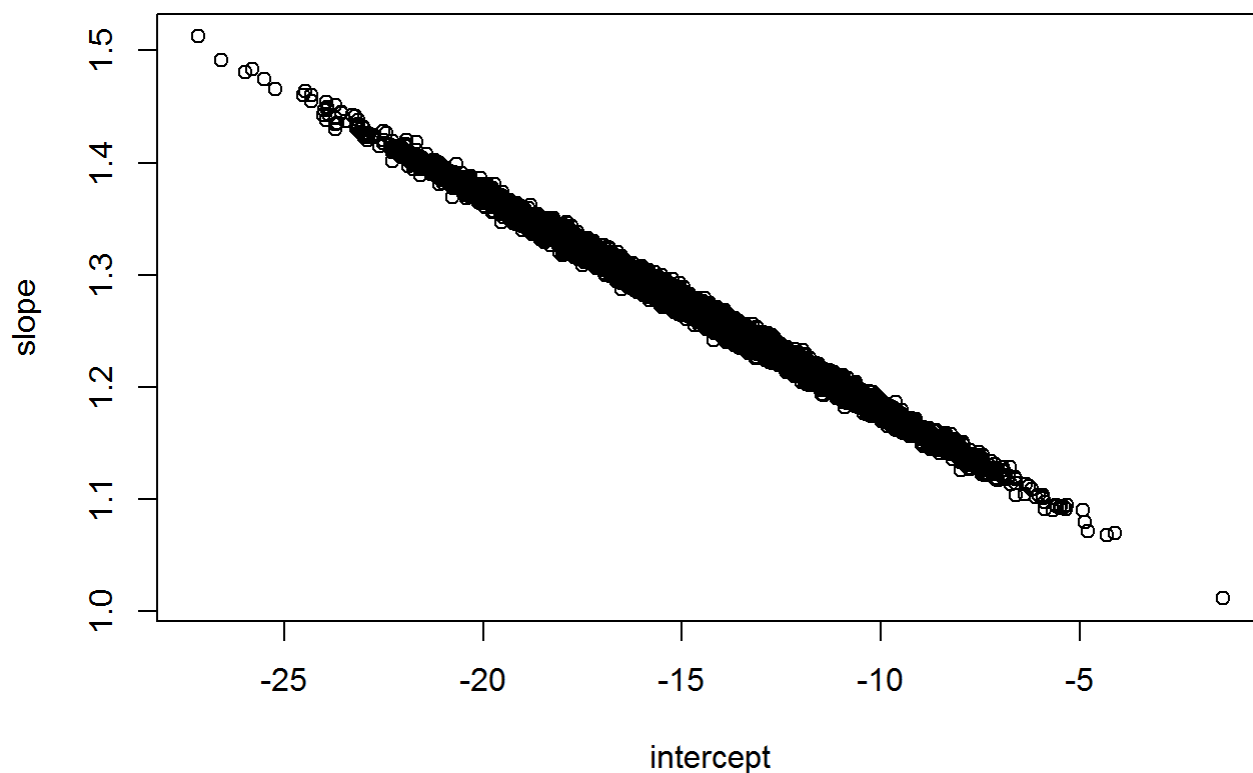
```
lm_model = lm(test$testA~test$testB)
est_intercept = coef(lm_model)[1]
est_slope = coef(lm_model)[2]
est_intercept
```

```
## (Intercept)
##      -14.4045
```

```
est_slope
```

```
## test$testB
##      1.263866
```

```
#Scatter Plot
plot(intercept, slope)
```



What does the plot of the estimated bootstrap slopes versus the estimated bootstrap intercepts plot indicate about these estimates?

Answer: The 10,000 points for the 10,000 samples roughly lie around a straight line. This shows that the estimates are fairly robust.

-1 point The slope is negative indicating that as the intercept goes up the slope goes down and vice versa.

12 points