

Enrichment Project #1: Rank-based Methods

*Yuan Gao, Kevin Lee, Akshay Govindaraj
Yijun (Emma) Wan, Peter Williams, Ruixuan Zhang
ygao390, kylee20, ywan40, agovindaraj6, pwilliams60, rzhang438 | @gatech.edu*

18 September 2018

Contents

1) Two-sample Studies (40%):	1
Bootstrap Re-sampling	2
Signed Rank Tests	3
Fligner Test for Equal Variances	6
Parametric F-test For Equal Variances	7
Parametric Two-Sample T-test	8
Goodness of Fit Tests	9
2) Multiple-Sample (ANOVA) Studies (60%):	11
Kruskal-Wallis Test	11
Friedman Test	11
Variance Testing	11
Parametric Two Sample Testing	11

1) Two-sample Studies (40%):

Locate a data set in the field of your interest, e.g., eCommerce, medical study, drug development, supply-chain/logistics operations, for applying the following procedures for two-sample studies.

For our two-sample study, we have selected a dataset published by the City of Chicago Transit Authority (CTA) that details daily ridership numbers (system-wide boarding totals) for both bus and rail services. As a pre-processing step, we take a subset of the original dataset that consists of weekday ridership only (weekends, and holidays excluded) for 2017 and 2018. Below is a quick preview of the dataset and pre-processing steps:

service_date	day_type	bus	rail_boardings	total_rides
2017-02-01	W	574,425	588,036	1,162,461
2017-02-02	W	865,433	695,663	1,561,096
2017-02-03	W	905,953	910,906	1,816,859
2017-02-06	W	740,684	373,851	1,114,535
2017-02-07	W	847,219	911,476	1,758,695

Note: For the purpose of comparison, we transformed the data to normalize the distributions for equal variances between them. This will allow us to perform both nonparametric and parametric tests to determine whether the two sample distributions are different. The variances of the samples were unequal, so we cannot nonparametric tests such as the Wilcoxon Rank Sum Test. With our new data, the variances are close enough that the distributions can then be compared.

Bootstrap Re-sampling

1. Calculate Pearson and Spearman coefficient of correlation and Kendall's Tau. Use a Bootstrap resampling procedure with $B = \text{\#bootstrap-samples} = 1000$ to assess the standard deviation (sd) of three estimates. Comment on your findings.

To begin we compute all three statistics using methods in base R, reported here:

```
## [1] "Pearson's Correlation: 0.647"
## [1] "Spearman's Correlation: 0.604"
## [1] "Kendall's Tau: 0.433"
```

We then compute these statistics using $B = 1000$ bootstrap samples based on our original dataset, which consists of 318 observations. Code and details are below. As shown in the table of results, bootstrap sample standard deviations for all correlation statistics were similar, and were slightly higher for spearman's rank:

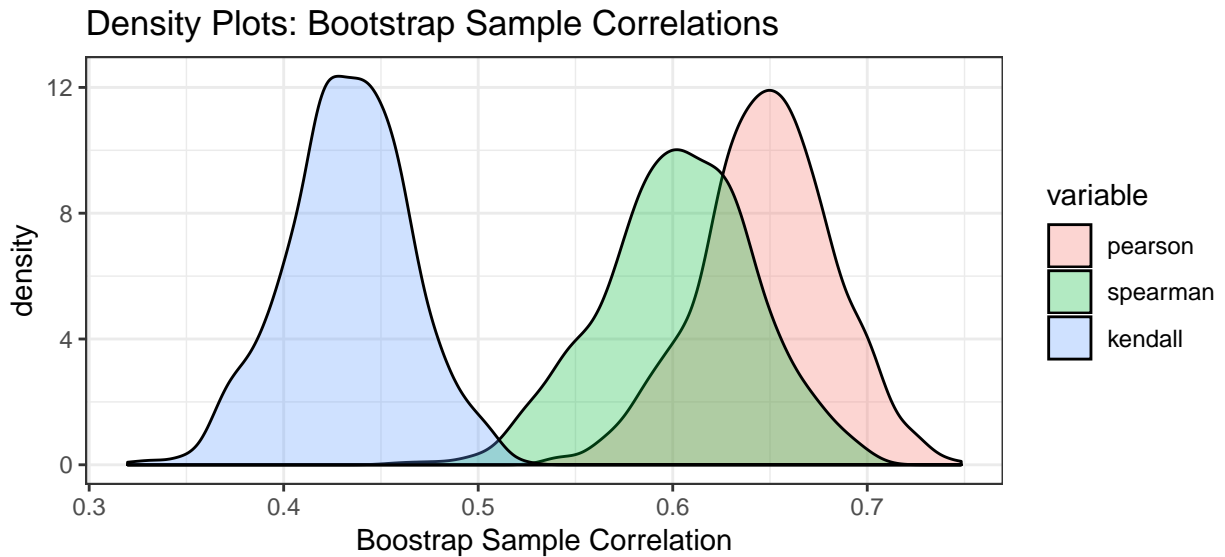
```
#set the number of bootstrap samples (with replacement - n = 318)
B <- 1000
#replicate process of collecting sample from original dataset, and compute sample stats
#results stored in: cor_result, with 1000 obs of bootstrap sample stats
#note this process takes a few seconds on a standard laptop
cor_result <- do.call('rbind',lapply(1:B, function(x){
  bootstrap_sample <- rider_data[sample(1:nrow(rider_data),
                                         size=nrow(rider_data), replace=T),]
  data.frame(pearson = cor(bootstrap_sample$bus,bootstrap_sample$rail_boardings,
                           method = 'pearson'),
             spearman = cor(bootstrap_sample$bus,bootstrap_sample$rail_boardings,
                           method = 'spearman'),
             kendall = cor(bootstrap_sample$bus,bootstrap_sample$rail_boardings,
                           method = 'kendall'))
}))
knitr::kable(round(data.frame(lapply(cor_result, sd)),digits=4),row.names=F,
             caption = 'Sample Standard Deviation: Bootstrap Sample Correlation')
```

Table 2: Sample Standard Deviation: Bootstrap Sample Correlation

pearson	spearman	kendall
0.0342	0.0391	0.0315

Density plots of these statistics are shown here, and provide visual evidence that the spearman rank correlation results and kendall's tau have more similar results than the parametric pearson's correlation:

```
melt_cor_result <- melt(cor_result, measure.vars = c('pearson','spearman','kendall'))
ggplot(melt_cor_result, aes(x = value, fill = variable)) + geom_density(alpha = 0.3) + theme_bw() +
  xlab('Bootstrap Sample Correlation') + ggtitle('Density Plots: Bootstrap Sample Correlations')
```



Signed Rank Tests

2. Apply Wilcoxon Signed Rank Test, Wilcoxon Sum Rank Test, Mann-Whitney U Test to compare two samples. For each test please state clearly what distribution is used to calculate the p-value.

Continuing on with the data given by the CTA on bus and rail ridership in 2017 and 2018, we can attempt to utilize the Wilcoxon Signed-Rank Test to determine whether there is a difference between paired data, but certain assumptions must first be met. In using the Wilcoxon Signed Rank Test, the data must be paired, the pairs must be independent, and the paired differences must be distributed symmetrically about 0 (the assumption that the pairs are drawn from the same population). The samples would also be better if from a nonnormal distribution (comparable parametric tests are more reliable only if data are normally distributed).

The Signed Rank Test (for two samples) looks at the difference in distributions between paired values, and our data are not paired in this case. Instead, the better test to use in the case of our CTA data would be the Wilcoxon Ranked Sum Test. The goal of the Wilcoxon Ranked Sum Test is to effectively compare two populations with a continuous response variable and nonnormal distributions. Thus, the Wilcoxon Ranked Sum Test is essentially analogous to a nonparametric form of a two sample t-test. Assuming the ridership count between the bus and rail systems are independent of each other and random, we can compare the two samples using the Wilcoxon Ranked Sum Test even without the assumption of paired data as was required in the Wilcoxon Signed Rank Test.

The Mann Whitney U Test accomplishes essentially the exact same purpose as the Wilcoxon Ranked Sum Test. The only differences are a differently labeled test statistic “U” compared to that of the Wilcoxon ranked sum test statistics “W” as well as the assumption that the shapes of the two population distributions have the same “shape”. They are even performed using the exact same function, with the “U” test statistic being equivalent to the “W” test statistics of the Wilcoxon Ranked Sum Test.

If we plot the density functions and histograms, we can see the distributions seem to resemble normal distributions, and they do resemble each other in terms of their “shape”.

```
#Performing the Wilcoxon Rank Sum Test/Mann Whitney U Test
with(rider_data, wilcox.test(bus,rail_boardings, paired=F))
```

```
##
## Wilcoxon rank sum test with continuity correction
##
## data: bus and rail_boardings
```

```
## W = 77742, p-value < 2.2e-16
## alternative hypothesis: true location shift is not equal to 0
```

The results of the Wilcoxon Ranked Sum Test show the p-value of the test to be small enough to be significant, meaning we reject the null hypothesis that the distributions of the two samples are identical. The fact that the Wilcoxon Ranked Sum Test and the Mann Whitney U Test are essentially the same test allows us to also state that the medians of the two distributions are also different.

For the calculation of the p-value, one must first be able to determine the test statistic being used in each test. For the Wilcoxon Rank Sum test, we find the

$$W_n = \sum_{i=1}^n iS_i(X, Y)$$

where n_1 is the sample size of the first sample (observations in bus ridership) , n_2 is the sample size of the second sample (observations in rail boarding), and $n_1 + n_2 = n$

S_i is an indicator function with value is 1 if ith ranked observation is from the first sample or 0 if from the second sample. [cite textbook]

The expected value of the W_n statistic is where the distribution will be centered around, and the variance of the statistic will be the symmetrical spread around that center:

$$E(W_n) = [n_1(n+1)]/2$$

$$Var(W_n) = [n_1n_2(n+1)]/12$$

As the test statistic is a linear rank statistic, we can then deduce that the W statistic is distributed approximately normally as long as the sample sizes are large enough (at least 10 observations per sample). We can then calculate the right-tail probability using Z-score with normal approximation with the following distribution:

$$W_n \sim N([n_1(n+1)]/2, [n_1n_2(n+1)]/12)$$

The Mann-Whitney U statistic is used to calculate p-value much in the same manner of the Wilcoxon Rank Sum Test. The test statistic can be calculated as such:

$$U = \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} D_{ij}$$

The test statistic U and W end up becoming the same value, used in the same approximately normal distribution as the W_n statistic, which means we can use the same tail probability to find the p-value.

In order to apply the Wilcoxon Signed-Rank Test, we will need to use a different set of data to demonstrate the test in R. For this purpose, we will use the data from this site: https://www.sheffield.ac.uk/polopoly_fs/1.569449!/file/stcp-Rdataset-Video.csv. A professor at the University of Sheffield collected data using “Likert” style questions to determine which of three new videos are most effective in informing the public of medical conditions. The four videos in questions are deemed the following: A, a new general video; B, a new medical profession video; C, the old video; and D, a demonstration using props. In this dataset, there are two particular variables of interest, and they are “TotalAGen” and “TotalDDemo”. These two variables are essentially the overall summed “Likert Scores” for each of the different types of videos. Our paired data will be between the scaled overall scores between video A’s reception and video D’s reception.

```

#For performing the Wilcoxon Signed Rank Test, download the dataset
video_data=read.csv('https://www.sheffield.ac.uk/polopoly_fs/1.569449!/file/stcp-Rdataset-Video.csv')
#Wilcoxon Signed Rank Test, with paired data
wilcox.test(video_data$TotalAGen,video_data$TotalDDemo, paired=T)

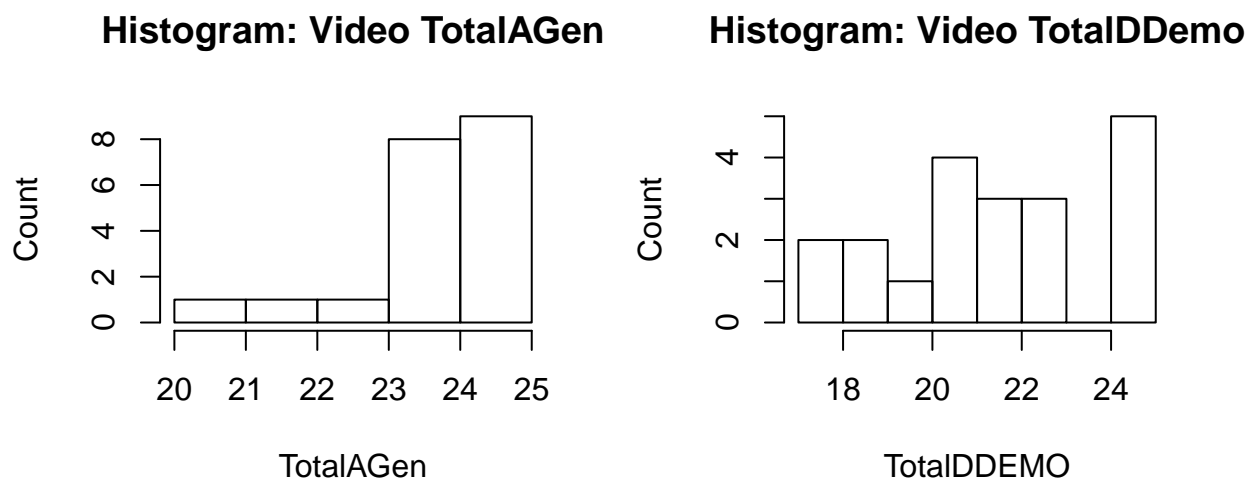
## Warning in wilcox.test.default(video_data$TotalAGen,
## video_data$TotalDDemo, : cannot compute exact p-value with ties

## Warning in wilcox.test.default(video_data$TotalAGen,
## video_data$TotalDDemo, : cannot compute exact p-value with zeroes

##
## Wilcoxon signed rank test with continuity correction
##
## data: video_data$TotalAGen and video_data$TotalDDemo
## V = 167.5, p-value = 0.003558
## alternative hypothesis: true location shift is not equal to 0

#Showing nonnormal distributions of the sample data
par(mfrow = c(1,2))
hist(video_data$TotalAGen, xlab = "TotalAGen", ylab = "Count", main = "Histogram: Video TotalAGen", bt,
hist(video_data$TotalDDemo, xlab = "TotalDDemo", ylab = "Count",main = "Histogram: Video TotalDDemo", b

```



We can assume the independence of each observation (person) as well as the fact that they should come from the same population (differences should be symmetric about 0). The data are also paired, so we can follow through with the test now that the assumptions are met. Through the Wilcoxon Signed Rank Test, we can see that the test statistic $V = 167.5$, which leads us to a p-value below our assumed alpha level of 0.05. Thus, we can reject the null hypothesis that the two sample distributions (Scores of Video A compared to scores of Video D) are the same.

The p-value of the signed rank test is based on a W test statistic and tabulated critical values when dealing with lower sample sizes, such that it can be calculated as follows:

$$W = \sum_{i=1}^n \text{sign}(x_{2,i} - x_{1,i}, |R_i|)$$

Where R_i is the rank of the observation and the statistic is the formula for the sum of the signed ranks. With this test statistic, one can determine a distribution with

$$E(W) = 0$$

$$Var(T) = [n(n+1)(2n+1)/6]$$

where n is the number of non-tied ranked observations. With this information, one can look at a table of set critical values dependent on sample size and quantiles with which to determine whether or not to reject the hypothesis that the

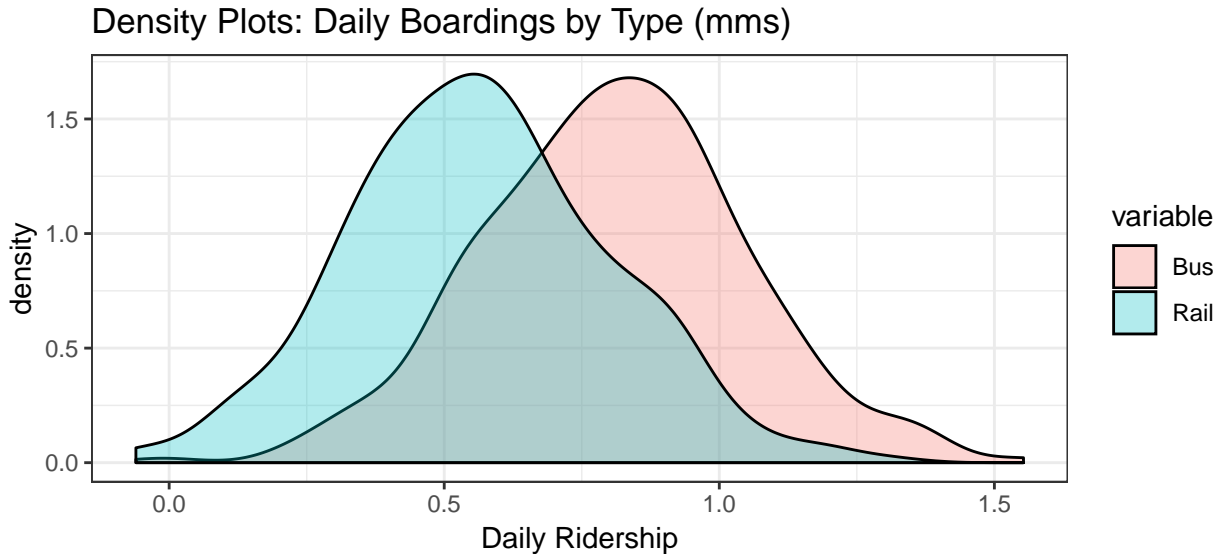
With larger sample sizes though, the distribution of the W statistic for the signed rank test also tends to follow normal approximation, which allows the use of z-score approximation to determine p-value through tail probability.

Fligner Test for Equal Variances

3. Use Fligner test for equal variances in these two samples. Explain how to calculate its p-value.

We again show a density plot that visualizes the counts of daily bus and rail ridership illustrating that variance of daily ridership seems homogenous across ride modality. Both ridership types also seems to be spread symmetrically around their central point as shown here:

```
boardings <- c(rider_data$bus, rider_data$rail_boardings)
factor_level <- factor(c(rep('bus',nrow(rider_data)), rep('rail',nrow(rider_data))))
ggplot(melted_rider_data, aes(x = value/1e06, fill = variable)) +
  geom_density(alpha = 0.3) + theme_bw() +
  xlab('Daily Ridership') + ggtitle('Density Plots: Daily Boardings by Type (mms)')
```



To explore this further we utilize the Fligner test to test the hypothesis $H_0 : \sigma_X^2 = \sigma_Y^2$, where samples X_1, \dots, X_{n_1} , and Y_1, \dots, Y_{n_2} refer to samples from daily ridership of bus and rail. The output of the test utilizing the 'fligner.test' function from R (package::car) is shown below, where our chosen $\alpha = 0.05$:

```
##
## Fligner-Killeen test of homogeneity of variances
##
## data:  boardings and factor_level
## Fligner-Killeen:med chi-squared = 0.0021818, df = 1, p-value =
## 0.9627
```

Since our test statistic $\chi_K^2 = 0.002$, is less than (3.8414588), we fail to reject our null hypothesis, that is $H_0 : \sigma_X^2 = \sigma_Y^2$.

The calculation of p-value is based on our realized χ^2 statistic which is computed:

$$\chi_{K^*}^2 = \frac{\sum_{j=1}^k n_j (\bar{a}_j - \bar{a})^2}{V^2}$$

Where k = number of groups, n_j is the number of observations for the j^{th} group, \bar{a}_j is the mean of the median centered, ranked, and subsequently normalized observations for the j^{th} group. \bar{a} is the mean of all median centered, ranked and normalized observations, and V^2 is the sample variance of the same normalized observations. If the assumptions are met, the distribution of this test statistic follows approximately the Chi-squared distribution with degrees of freedom $k - 1$. We can use the Chi-square distribution to get the p-value of Fligner test given our realized sample statistic.

Parametric F-test For Equal Variances

4. Use the parametric F-test for equal variances to the data; comment on the difference of the assumptions and results compared to them in (iii).

F-test for testing equality of variance is used to test the hypothesis of the equality of two population variances. The test statistic can be obtained by computing the ratio of the two variances S_A^2 and S_B^2 .

$$F = \frac{S_A^2}{S_B^2}$$

The degrees of freedom are $n_A - 1$ (for the numerator) and $n_B - 1$ (for the denominator). And, the more this ratio deviates from 1, the stronger the evidence for unequal population variances.

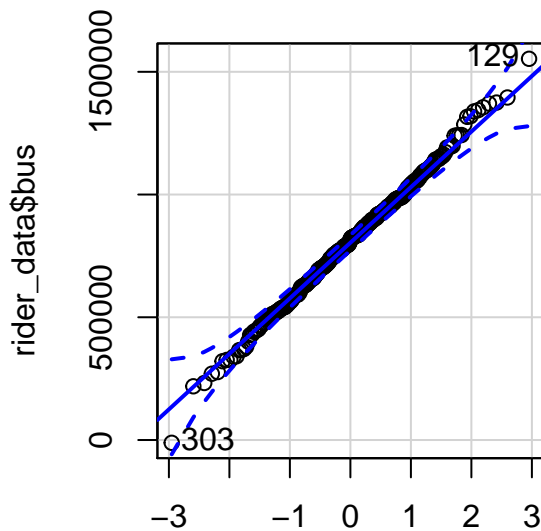
```
##
## F test to compare two variances
##
## data: rider_data$bus and rider_data$rail_boardings
## F = 1.0082, num df = 317, denom df = 317, p-value = 0.9422
## alternative hypothesis: true ratio of variances is not equal to 1
## 95 percent confidence interval:
## 0.8086292 1.2569753
## sample estimates:
## ratio of variances
## 1.00818
```

The result shows that the p-value = 0.01772 which is smaller than our $\alpha = 0.05$. In conclusion, there is significant difference between the two variances.

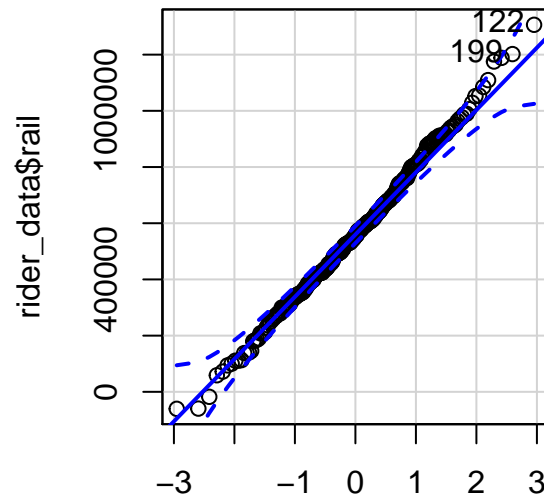
Assumptions for Conover and F-test: According to the above, the Conover test does not assume that all populations are normally distributed. However, F-test is very sensitive to departure from the normal assumption. We use Q-Q plot (quantile-quantile plot) to graphically evaluate the normality of a variable.

```
## [1] 303 129
```

Q-Q Plot – Bus Boardings



Q-Q Plot – Rail Boardings



```
## [1] 122 199
```

As we can see, there are many points don't fall approximately along this reference line, for both rail and bus boardings, so, we cannot assume normality for either sample.

Results for Conover and F-test: Although the results are consistent based on this sample, Conover test is recommended when the normality assumption is not viable. As a result, use F-test on the data has shortcomings compared to use Conover test.

Parametric Two-Sample T-test

5. Depending on the outcomes from the F-test in (iv), apply an appropriate parametric two- sample t-test to the data; comment on the difference of the assumptions and results compared to them in (ii).

The aim of the two-sample t-test is to find the difference between the two sample means in comparing their respective populations. In utilizing the two sample t-test, some assumptions must first be met. One must assume that the observations of our data are independent of each other, similar to the nonparametric tests. Another such assumption is for the response variable to be continuous, approximately normal distribution. Our normalized ridership data allows us to be able to perform the parametric (pooled) two sample t-test. Similar to the nonparametric tests performed (Mann-Whitney U), the data must be continuous.

```
##
## Two Sample t-test
##
## data: rider_data$bus and rider_data$rail_boardings
## t = 12.917, df = 634, p-value < 0.00000000000000022
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## 206160.5 280083.7
## sample estimates:
## mean of x mean of y
## 807865.8 564743.7
```


As we can see from the results of the test, the p-value is quite low, showing that there is a significant difference between the two means. In comparing the two tests, the parametric t-test tests for significance in difference between the means of the two samples, while the nonparametric tests (Mann-Whitney U) test for significance in differences between the distributions (and medians) of the two samples. In terms of results, both the Mann-Whitney U test and the two sample t-test show a significant difference in mean and median for the two samples. Despite the fact that they both could determine whether there are differences between the two sample distributions, the two sample t-test is a more reliable test if only for the fact that the Mann-Whitney U test is better suited for nonnormal distributions (while we transformed our data to become more normal so as to meet the assumption of normality for the t-test). The t-test is a more reliable test compared to its nonparametric analog when the distributions are approximately normal.

Goodness of Fit Tests

6. Apply Kolmogorov-Smirnov, Anderson-Darling, Cramer-Von Mises, Shapiro-Wilk Tests for normality to the two samples separately; comment on the findings by comparing results obtained from these four tests. Make a statement about the situation that a particular procedure might be more appropriate. Moreover, based on the results learned here, comment on whether the parametric methods used in (iv) and (v) are appropriate.

Kolmogorov-Smirnov:

```
ks.test(rider_data$bus,"dnorm",mean(rider_data$bus),sd(rider_data$bus))

##
## One-sample Kolmogorov-Smirnov test
##
## data: rider_data$bus
## D = 1, p-value < 0.00000000000000022
## alternative hypothesis: two-sided

ks.test(rider_data$rail_boardings,"dnorm",mean(rider_data$rail_boardings),sd(rider_data$rail_boardings))

##
## One-sample Kolmogorov-Smirnov test
##
## data: rider_data$rail_boardings
## D = 1, p-value < 0.00000000000000022
## alternative hypothesis: two-sided
```

Anderson-Darling:

```
ad.test(rider_data$bus)

##
## Anderson-Darling normality test
##
## data: rider_data$bus
## A = 0.16818, p-value = 0.9357

ad.test(rider_data$rail_boardings)

##
## Anderson-Darling normality test
##
## data: rider_data$rail_boardings
## A = 0.32229, p-value = 0.5267
```

Cramer-Von Mises:

```
cvm.test(rider_data$bus)
```

```
##  
## Cramer-von Mises normality test  
##  
## data: rider_data$bus  
## W = 0.023255, p-value = 0.9324
```

```
cvm.test(rider_data$rail_boardings)
```

```
##  
## Cramer-von Mises normality test  
##  
## data: rider_data$rail_boardings  
## W = 0.054268, p-value = 0.4491
```

Shapiro-Wilk:

```
shapiro.test(rider_data$bus)
```

```
##  
## Shapiro-Wilk normality test  
##  
## data: rider_data$bus  
## W = 0.99794, p-value = 0.965
```

```
shapiro.test(rider_data$rail_boardings)
```

```
##  
## Shapiro-Wilk normality test  
##  
## data: rider_data$rail_boardings  
## W = 0.99658, p-value = 0.7343
```

Comparison of results: All the above methods (Kolmogorov-Smirnov, Anderson-Darling, Cramer-Von Mises, Shapiro-Wilk Tests) have the same result: reject H_0 under our sample. But the p-value of each test is significantly different which means the measurement for each one is different.

Appropriate situation for the four tests: Shapiro-Wilk test is the most powerful test for all types of distribution whereas Kolmogorov-Smirnov test is the least powerful test. The performance of Anderson-Darling test is quite comparable with Shapiro-Wilk test. Cramer-Von Mises test is an alternative to the Kolmogorov-Smirnov test. However, the power of Shapiro-Wilk test is still low for small sample size.

Kolmogorov-Smirnov is based on the empirical distribution function (ECDF), and the maximum distance between these two curves. So, it is independent with the underlying cumulative distribution function being tested. But, it is sensitive on the center of the distribution than at the tails. It is suitable for small samples, ties are no problem and has omnibus test, but it is low power if prerequisites are not met. The Kolmogorov-Smirnov test statistic is defined as

$$D = \max_{1 \leq i \leq N} (F(Y_i) - \frac{i}{N}, \frac{i}{N} - F(Y_i))$$

Anderson-Darling test is used to test samples with a specific distribution. It is a modification of the Kolmogorov-Smirnov (K-S) test and gives more weight to the tails than does the K-S test. It is high power when testing for normal distribution but is statistic based on squares. The Anderson-Darling test statistic is defined as

$$A^2 = -N - S$$

where

$$S = \sum_{i=1}^N \frac{(2i-1)}{N} [\ln F(Y_i) + \ln (1 - F(Y_{N+1-i}))]$$

Shapiro-Wilk test calculates a W statistic that tests whether a random sample, x_1, x_2, \dots, x_n comes from (specifically) a normal distribution. It is highest power among all tests for normality but test for normality only. The W statistic is calculated as follows:

$$W = \frac{(\sum_{i=1}^n a_i x_{(i)})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Cramer-Von Mises is higher power than KS test, but it's a statistic based on squares. Cramer-Von Mises statistic is defined as:

$$U^2 = T - n(\bar{F} - \frac{1}{2})^2$$

where

$$\bar{F} = \frac{1}{n} \sum F(x_i)$$

As a result, due to our data is not normally distributed and we have more than 300 data points, so in our opinion, it is more reasonable to choose Anderson-Darling test.

2) Multiple-Sample (ANOVA) Studies (60%):

Locate one data set each for the two problems below in the field of your interest, e.g., eCommerce, medical study, drug development, supply-chain/logistics operations, for applying the following procedures for ANOVA studies.

Kruskal-Wallis Test

1. Apply Kruskal-Wallis Test for an one-way ANOVA study. If it is suitable, perform a K- W pairwise comparisons. Make conclusions about your findings.

Friedman Test

2. Use Friedman test and also the F-Test discussed in the textbook page 148 for the study of one-way ANOVA with one blocking variable. Comment on your findings. If it is suitable, perform a K-W pairwise comparisons. Make conclusions about your findings.

Variance Testing

3. Conduct a variance test based on the procedure (Conover test) given in Section 8.3 textbook. Comment on your findings.

Parametric Two Sample Testing

4. Repeat the same studies in (i), (ii) and (iii) using parametric approaches (also include the possible pairwise comparisons). State the assumptions needed for the parametric approaches. Compare the results here against those in (i), (ii) and (iii), respectively. Note that if there are certain assumptions (e.g., normality and equal-variance) required in the parametric studies, please apply appropriate procedures to “test” the assumption.