

UNIT 4 HW

This class allows you to practice preparing professional looking reports. Make sure all reports are typed and all graphs (unless otherwise noted) are computer generated and copied and pasted into your report. If you would like help with Word or Excel please don't hesitate to ask.

1. Read Chapter 4 from Statistical Sleuth and answer the conceptual problems at the end of the chapter. Note: You do not need to type these up and turn them in. The answers are at the very end of the chapter.
2. When wildfires ravage forests, the timber industry argues that logging the burned trees enhances forest recovery; the EPA argues the opposite. The 2002 Biscuit Fire in southwest Oregon provided a test case. Researchers selected 16 fire-affected plots in 2004, before any logging was done and counted tree seedlings along a randomly located transect pattern in each plot. They returned in 2005, after nine of the plots had been logged, and counted the tree seedlings along the same transects. The percent of seedlings lost from 2004 to 2005 is recorded in the table below for logged (L) and unlogged (U) plots:

Test the EPA's assertion (and thus the opposite of the logging industries assertion) that logging actually increases the percentage of seedlings lost from 2004 to 2005.

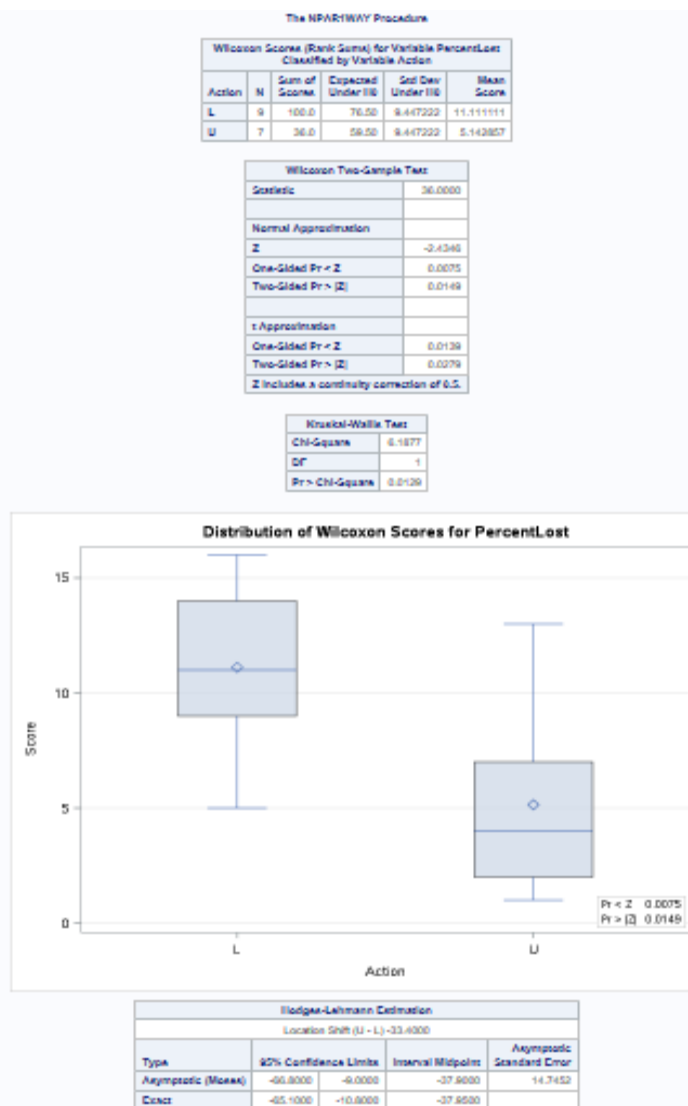
- a. Perform a complete analysis using a rank sum test in SAS. (Logging data).

```
proc import datafile = '/home/chec0/New Folder/Logging.csv'
  out = logging
  dbms = CSV
;

proc print data = logging;
run;

proc npar1way data=logging wilcoxon;
var PercentLost;
class Action;
exact HL;
run;
```

CV
-1.64592439



Step 1:

H_0 : The distribution of the percent of seedlings in the logged plots is equal to that of the unlogged plots

H_a : The distribution of the percent of seedlings in the logged plots is less than that of the unlogged plots.

Step 2:

$$\bar{R} = 8.5 \quad s_R = 30.2812$$

Mean(T) for L group = $9 * 8.5 = 76.5$ & Mean(T) for U group = $7 * 8.5 = 59.5$

SD(T) = $Z = 9.4472$

Step3: Z-test = -2.4346

Step4: P-value = .0075(1-sided) & .0149(2-sided)

Step5: Reject H_0

Step6: There is sufficient evidence to suggest that the distribution of the tree seedlings in the logged plots is less than that of the unlogged plots. (the p-value=.0075 from a one sided rank sum test). A 95% confidence interval based off rank

sum test for the increase in median percent seedlings in favor of the unlogged plots is (10.8%, 65.1%)

Scope of inference: Since the plots were not randomized to receive either the logging or not logging treatment, no causation can be implied here. Since the transect patterns were randomly selected, this inference can be generalized to only the 16 plots.

- b. Verify the p-value and confidence interval by running the rank sum test in R (using R function `Wilcox.test`). (You do not need to repeat the complete analysis ... simply cut and paste a screen shot of your code and the output.) You may use: <https://www.r-bloggers.com/wilcoxon-mann-whitney-rank-sum-test-or-test-u/> for reference.

```
> L = c(45, 53.1, 40.8, 75.5, 46.7, 85.4, 85.6, 18.2)
> U = c(23.6, 13.3, 34.2, 18.1, 56.1, -8.1, -20.1)
> wilcox.test(L,U, correct=FALSE)
```

Wilcoxon rank sum test

data: L and U
W = 49, p-value = 0.01399
alternative hypothesis: true location shift is not equal to 0

3. Conduct a Welch's two-sample t-test on the Education Data from HW 3 (untransformed). Perform a complete analysis using SAS to test the claim that the mean income of college educated people (16 years of education) is greater than the mean of those with a high school education only (12 years of education).

- a. **State the problem, address the assumptions. Be sure to support with your knowledge of theory (CLT) as well as with histograms, box plots, q-q plots, etc.**

Problem: Test the claim that the distribution of incomes for those with 16 years of education(μ_1), exceeds the distribution for those with 12 years of education(μ_2).

Assumptions: In order to test the claim that the "distribution" of incomes of those with 16 years of education exceed the "distribution" of incomes with 12 years of education we will focus on the location parameters: the median. We will test if there is sufficient evidence to suggest that the median income of those with 16 years of education exceeds the median income of those with only 12 years of education. The histogram and box plot below indicate strong evidence of inequality of variance between the two populations. With a large sample size the assumption will be robust and will assume independence and run Welch's t test.

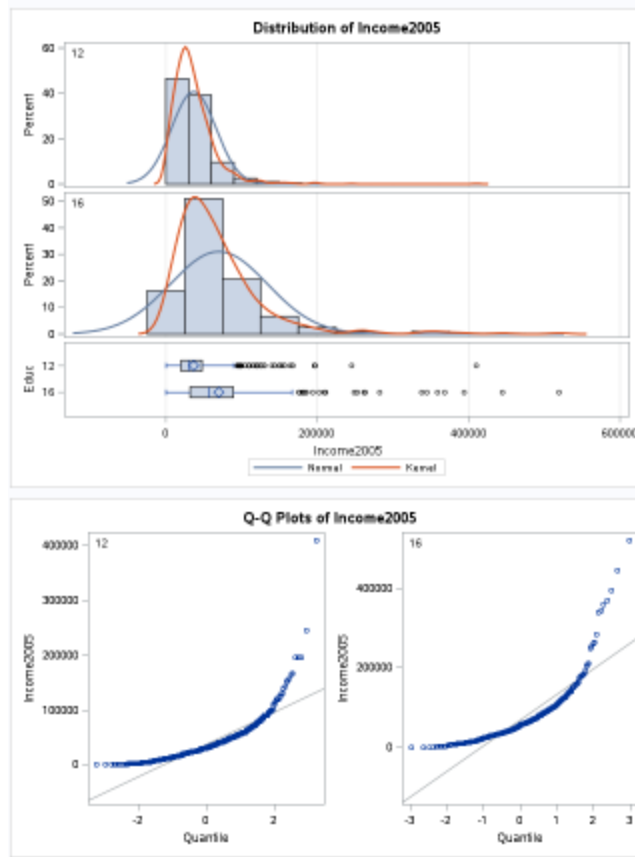
The TTEST Procedure
Variable: Income2005

Educ	Method	N	Mean	Std Dev	Std Err	Minimum	Maximum
12		1020	30864.9	29369.7	919.6	300.0	413006
16		406	69967.0	64256.8	3189.0	200.0	518340
Diff (1-2)	Pooled		-33132.1	42326.9	2483.8		
Diff (1-2)	Satterthwaite		-33132.1		3219.0		

Educ	Method	Mean	95% CL Mean	Std Dev	95% CL Std Dev
12		30864.9	30864.9	29369.7	26148.2
16		69967.0	69967.0	64256.8	60120.1
Diff (1-2)	Pooled	-33132.1	-33132.1	42326.9	40628.0
Diff (1-2)	Satterthwaite	-33132.1	-33132.1	3219.0	43810.9

Method	Variance	DF	t Value	Pr > t
Pooled	Equal	1424	-13.34	<.0001
Satterthwaite	Unequal	473.85	-9.98	<.0001

Equality of Variances				
Method	Num DF	Den DF	F Value	Pr > F
Folded F	405	1019	4.79	<.0001



b. Show all 6 steps, including a thoughtful, thorough, yet non-technical conclusion. Include a confidence interval.

Step1: $H_0: \mu_1 = \mu_2$ $H_a: \mu_1 > \mu_2$

Step2: CV = -1.645

Step3: Value of Test Statistic (Satterwaite): $t = -9.98$

Step4: p-value : $p < 0.0001$

Step5: Reject H_0

Step6: There is overwhelming evidence at the alpha = 0.05 level of significance from sample 1 sided Welch's t test ($p < 0.0001$) that the mean income in 2005 for people with 12 years of education is less than mean income for those that had 16 years of

education. A 95% confidence interval for the difference is [\$26610, \$39653] in favor of those with college degrees.

- c. **Include a scope of inference at the end. (You may copy and paste this from a previous HW if you like.)**

This was an observational study, and thus we cannot confirm that the years of education caused the increase in income, only that they are associated with each other. There is little detail about the randomness of the sample, although it is doubtful that it was a random sample. We must limit the inference gained from this study to only the subjects of this sample.

- d. **Verify the Welch's t statistic and p-value with R (using R function t.test). Simply cut and paste your R code and output. You may use:**

http://rcompanion.org/rcompanion/d_02.html for reference.

```
data <- read.csv(file="C:/Users/che/Desktop/R/EducationData.csv", header=TRUE, sep=",")
t.test(Income2005 ~ Educ, data=data,
       var.equal=FALSE,
       conf.level=0.95)
```

Welch Two Sample t-test

data: Income2005 by Educ

t = -9.9827, df = 473.85, p-value < 2.2e-16

alternative hypothesis: true difference in means is not equal to 0

95 percent confidence interval:

-39653.77 -26610.39

sample estimates:

mean in group 12 mean in group 16

36864.90 69996.97

- e. **Would you prefer to run the log transformed analysis you ran in HW3, or do you feel this analysis is more appropriate? Why or Why not? (Make mention of the assumptions as well as the parameters that each test provides inference on. As you know, they are different.)**

The original distributions of income were right-skewed, the median may be a more valuable measure of center. So for this reason a log transformation with this test analysis may be preferred.

4.

- a. Chapter 4, Problem 20 from the text. Show all work. "By hand" here means actually by hand. Simply take a picture of your work and include it in your pdf/doc file. Include your sorted, labeled, and ranked data; your calculations of the mean and standard deviation of the assumed distribution of the rank sum statistic under H_0 ; your calculation of the Z

statistic with a continuity correction; your p-value, and conclusion. (No confidence interval necessary here.)

	GROUP	ORDER	RANK
18.8	NT	1	1
20	NT	2	2
20.1	NT	3	3
20.9	NT	4	4.5
20.9	NT	5	4.5
21.4	NT	6	6
22	T	7	7
22.7	NT	8	8
22.9	NT	9	9
23	T	10	10
24.5	T	11	11
25.8	T	12	12
30	T	13	13
37.6	T	14	14
38.5	T	15	15

$T \text{ GROUP 1} = 7 + 10 + 11 + 12 + 13 + 14 + 15 = 82$
 $NT \text{ GROUP 2} = 1 + 2 + 3 + 4.5 + 4.5 + 6 + 8 + 9 = 38$

$R = 8 \quad SR = 4.46814$
 $MEAN(T) \text{ GROUP 1} = 7 \times 8 = 56$
 $\text{GROUP 2} = 8 \times 8 = 64$

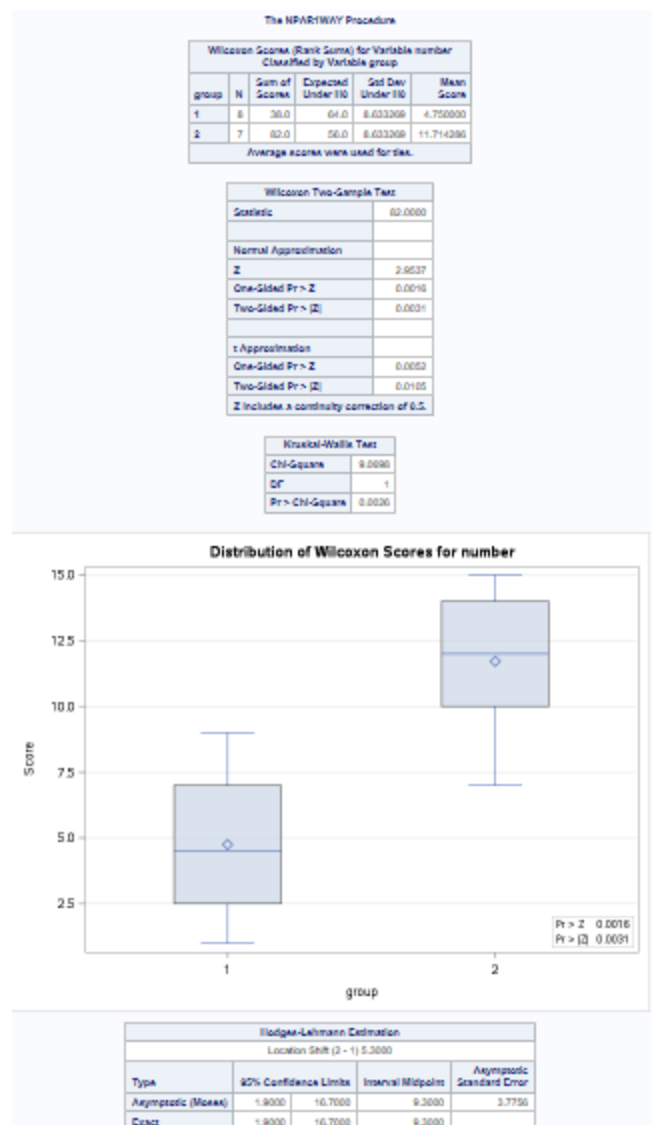
$SD(T) = 4.46814 \sqrt{\frac{7 \times 8}{7 + 8}} = 8.63327$
 $P\text{-VALUE} = .0013$

$Z = \frac{(38.5 - 64)}{8.63327} = -2.9537$
 $\frac{(81.5 - 56)}{8.63327} = 2.9537$

- b. Problem 21 from the text. Take a screen capture of the SAS output in addition to your response.

Sum of Ranks SD(T) and Z statistic are the same

P-value is .0016 in SAS and .0013 by hand



- c. Write up a complete analysis using the information you have gained from A and B to test the claim that the distributions are different.

- i. State the problem.

Test the claim that the mean of non-trauma patients exceeds the mean trauma patients

- ii. State the assumptions you are making and why you are making them. Justify your decisions. Print out any histograms, q-q plots, box plots, etc. that you use in your justification.
- iii. Show all 6 steps of the hypothesis test for the rank sum test of the trauma data. Use the critical values, test statistics, p-values, etc. obtained above. Add a confidence interval from the Hodges-Lehmann procedure (from SAS).

Step1:

H_0 : The distribution of the "trauma(group1)" metabolic expenditures is the same as the distribution of the "non-trauma(group2)" metabolic expenditures.

H_a: The distribution of the "trauma(group1)" metabolic expenditures is different from the distribution of the "non-trauma(group2)" metabolic expenditures

Step2:

$$\bar{R} = 8 \quad s_R = 4.46814$$

$$\text{Mean}(T) \text{ for group1} = 7 \cdot 8 = 56 \text{ \& Mean}(T) \text{ for group2} = 8 \cdot 8.5 = 64$$

$$\text{SD}(T) = 4.46814 \sqrt{(7 \cdot 8)/(7+8)} = 8.63327$$

Step3:

$$\text{Z test for group2} = (38.5 - 64)/8.63327 = -2.9537$$

$$\text{for group1} = (81.5 - 56)/8.63327 = 2.9537$$

Step4: P-value = .0016(1-sided) & .0031(2-sided)

Step5: Reject the H₀

Step6: There is evidence at the alpha = 0.05 level of significance from Welch's t test($p = 0.0016$) that the mean metabolic expenditures of group2(Non-Trauma) is less than mean metabolic expenditures of group2(Trauma). A 95% confidence interval for the difference is [1.9000, 16.7000] in favor of those with Trauma

iv. Also include a scope of inference statement.

Since the metabolic expenditures were not randomized to receive either the trauma or non-trauma, no causation can be implied here. Since the metabolic expenditures were randomly selected, this inference can be generalized to only the 15 patients.

5. A study was performed to test a new treatment for autism in children. In order to test the new method, parents of children with autism were asked to volunteer for the study in which 9 parents volunteered their children for the study. The children were each asked to complete a 20 piece puzzle. The time it took to complete the task was recorded in seconds. The children then received a treatment (20 minutes of yoga) and were asked to complete a similar but different puzzle. The data from the study is below:

Child	Before	After
1	85	75
2	70	50
3	40	50
4	65	40
5	80	20
6	75	65
7	55	40
8	20	25
9	70	30

- a. Calculate the statistic S for a signed rank test by hand showing the final table with the absolute differences, the signs, and the ranks. Also, show your calculation of the z-statistic (standardized S statistic).

Child	BEFORE	AFTER	difference	SIGN	RANK
8	20	25	5	+	1
1	85	75	10	-	3
3	40	50	10	+	3
6	75	65	15	+	5
7	55	40	20	+	6
2	70	50	25	+	7
4	65	40	40	+	8
9	70	30	60	+	9
5	80	20			
Sum	OF	+RANKS	41		
"	"	-RANKS	4		

$$Z = \frac{41 - .5 - (9 \times 10) / 4}{\sqrt{\frac{9 \times 10 \times 19}{24}}} = 2.13$$

$$P(Z > 2.13) = 0.016586 = 0.0166$$

b. Verify your calculation in both SAS and R. Simply cut and paste your code and relevant output.

SAS CODE and Output

```
data autism;
```

```
input Child Before After;
```

```
datalines;
```

```
1 85 75
```

```
2 70 50
```

```
3 40 50
```

```
4 65 40
```

```
5 80 20
```

```
6 75 65
```

```
7 55 40
```

```
8 20 25
```

```
9 70 30
```

```
;
```

```
data autism2; set Child;
```

```
diff = Before - After;
```

```
run;
```

```
proc univariate data=autism2;
```

```
var diff;
```

```
run;
```

Tests for Location: Mu0=0				
Test	Statistic		p Value	
Student's t	t	2.540341	Pr > t	0.0347
Sign	M	2.5	Pr >= M	0.1797
Signed Rank	S	18.5	Pr >= S	0.0313

R Code and Output

```
before <- c(85, 70, 40, 65, 80, 75, 55, 20, 70)
```

```
after <- c(75, 50, 50, 40, 20, 65, 40, 25, 30)
```

```
##Sign Test
```

```
binom.test(6, 9, alternative='greater')
```

```
##Signed Rank Test
```

```
wilcox.test(before, after, paired=T, alternative='greater')
```

```
Wilcoxon signed rank test with continuity correction

data: before and after
V = 41, p-value = 0.01618
alternative hypothesis: true location shift is greater than 0
```

c. Conduct the six step hypothesis test using your calculations from above to test the claim that the yoga treatment was effective in reducing the time to finish the puzzle.

Step1:

H_0 : The **median** difference in time to finish a puzzle before and after yoga is zero

H_A : The **median** difference in the time to finish a puzzle after yoga is greater than zero (if one-sided).

Step2: Critical Value 1.96 one sided

Step3: zstat = 2.13

Step4: P-value =0.0166

Step5: **Reject the H_0**

Step6: **There is sufficient evidence to suggest that the median time to finish a puzzle for individuals before treatment is greater than the median time to finish a puzzle for individuals after they do yoga treatment.(p-value 0.0166 from one sided Signed Rank Test). There is evidence that yoga is associated with shorter puzzle solution times.**

d. Use SAS to conduct a six step hypothesis test using a paired t-test to test the claim that the yoga treatment was effective in reducing the time to finish the puzzle.

```
data autism;
```

```
input Child Before After;
```

```
datalines;
```

```
1 85 75
```

```
2 70 50
```

```
3 40 50
```

```
4 65 40
```

```
5 80 20
```

```
6 75 65
```

```
7 55 40
```

```
8 20 25
```

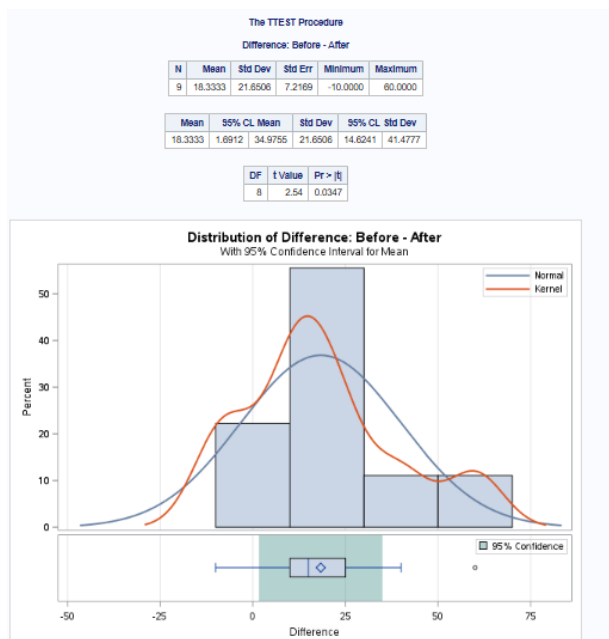
```
9 70 30
```

```
;
```

```
proc ttest data= autism;
```

```
paired Before*After;
```

```
run;
```



Step1:

H_0 : The **median** difference in time to finish a puzzle before and after yoga is zero

H_A : The **median** difference in the time to finish a puzzle before and after yoga is greater than zero (if one-sided).

Step2: Critical value 1.96 for one sided

Step3: zstat =2.13

Step4: take half of 0.0347 and p-value is .01735 <.05

Step5: Reject the H_0

Step6: There is sufficient evidence to suggest that the median time to finish a puzzle for individuals before treatment is greater than the median time to finish a puzzle for individuals after they do yoga.(p-value 0.01735 from one sided Paired T Test). There is evidence that yoga is associated with shorter puzzle solution times.

e. Verify your calculations in R. Simply cut and paste your code and relevant output.

Input = ("

Child Before After

1 85 75

2 70 50

3 40 50

4 65 40

5 80 20

6 75 65

7 55 40

8 20 25

9 70 30

")

```
Data = read.table(textConnection(Input),header=TRUE)
```

```
t.test(Data$Before,
```

```
      Data$After,
```

```
      paired=TRUE,
```

```
      conf.level=0.95)
```

```
      Paired t-test

data:  Data$Before and Data$After
t = 2.5403, df = 8, p-value = 0.03469
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 1.691182 34.975485
sample estimates:
mean of the differences
      18.33333
```

f. Use your data from above to construct a “complete analysis” of the test that you feel is most appropriate to test the claim that the yoga treatment was effective in reducing the time to finish the puzzle. This is simply formatting your results. You should be able to cut and paste most of the work from above.

Problem: Test the claim that the yoga treatment was effective in reducing the time for autism patients in finishing a puzzle.

Assumptions: Data are normally distributed, independent between observations with a small sample size. I will use the Signed Rank Test.

Step1:

H_0 : The **median** difference in time to finish a puzzle before and after yoga is zero

H_A : The **median** difference in the time to finish a puzzle after yoga is greater than zero (if one-sided).

Step2: Critical Value 1.96 one sided

Step3: $z_{stat} = 2.13$

Step4: P-value = 0.0166

Step5: Reject the H_0

Step6: There is sufficient evidence to suggest that the median time to finish a puzzle for individuals before treatment is greater than the median time to finish a puzzle for individuals after they do yoga treatment. (p-value 0.0166 from one sided Signed Rank Test). There is evidence that yoga is associated with shorter puzzle solution times.

BONUS (1 pt on 20 pt scale, 5pts on 100 point scale, etc.) This one is challenging and involves hard core SAS coding! Using our permutation test SAS code that we have used in prior HWs, do the following:

- Build the permutation distribution for the rank sum statistic for the Trauma data used above. Use 5000 permutations. Use SAS to fit / overlay a normal curve to the resulting histogram. Compare the mean and standard deviation of this normal curve that was fit to the permutation / randomization distribution to the mu and sigma you found in earlier in the homework.
- Compare the one-sided p-value found in this permutation distribution with the one found in prior questions.

HINT: Don't mind the highlight; the whole thing is the hint. You will need to work code similar to what is to the right into the permutation test SAS code we used before (in place of Proc ttest). You will also have to do some research on how to get your hands on the sum of the ranks statistic (a good start is to print the outnpar data set!).

```
ODS OUTPUT WilcoxonTest = outnpar;  
PROC NPARIWAY DATA=learn WILCOXON;  
  CLASS group;  
  VAR score;  
  EXACT;  
RUN;  
PROC PRINT DATA=outnpar;  
RUN;
```