

STA 104 Applied Nonparametric Statistics

Chapter 2: One-Sample Methods for Location Problem

Xiner Zhou

Department of Statistics, University of California, Davis

Table of contents

1. Wilcoxon Signed Rank Test
2. Signed Test
3. A Comparison of Statistical Tests
4. Paired Comparisons

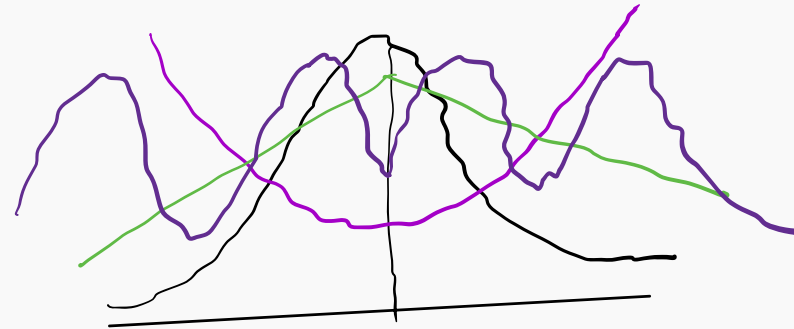
Wilcoxon Signed Rank Test

Setting

Suppose we have a random sample $x_1 \dots x_n$, i.e. data

- The x 's are mutually independent.
- they are from a population that is
continuous
symmetric about the median θ

Assumptions



Two-Sided Test:

$$H_0 : \theta = \theta_0 \text{ versus } H_a : \theta \neq \theta_0$$

One-Sided Upper-Tail Test:

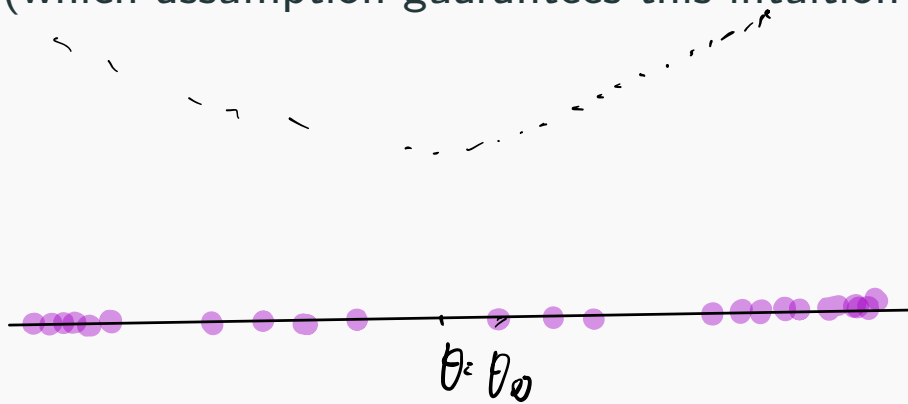
$$H_0 : \theta = \theta_0 \text{ versus } H_a : \theta > \theta_0$$

One-Sided Lower-Tail Test:

$$H_0 : \theta = \theta_0 \text{ versus } H_a : \theta < \theta_0$$

Motivation

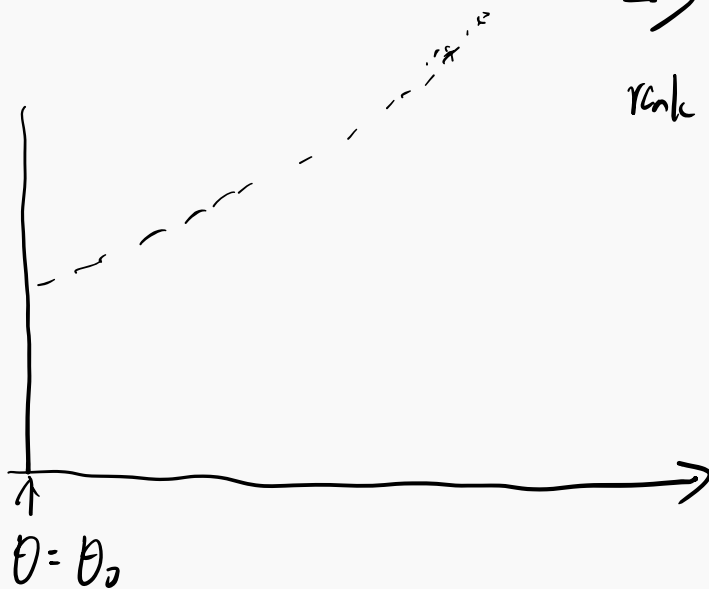
Intuition: If θ_0 was the true median of the population, then the magnitude in terms of absolute value of the centered data is nothing to do with the sign of the centered data (which assumption guarantees this intuition valid?)



Assume : Symmetri dist.

Motivation

- **Centering:** subtract θ_0 from each observation x_1, \dots, x_n to form a modified sample $x'_1 = x_1 - \theta_0, \dots, x'_n = x_n - \theta_0$ *centered data*
- **Flip:** form absolute values $|x'_1| \dots |x'_n|$
- **Rank:** Order them from least to greatest, let R_i denote the rank of i th observation $|x'_i|$



\Rightarrow

1	2	3	4
\uparrow	\uparrow	\uparrow	\uparrow
rank: 1	2	3	4

Motivation

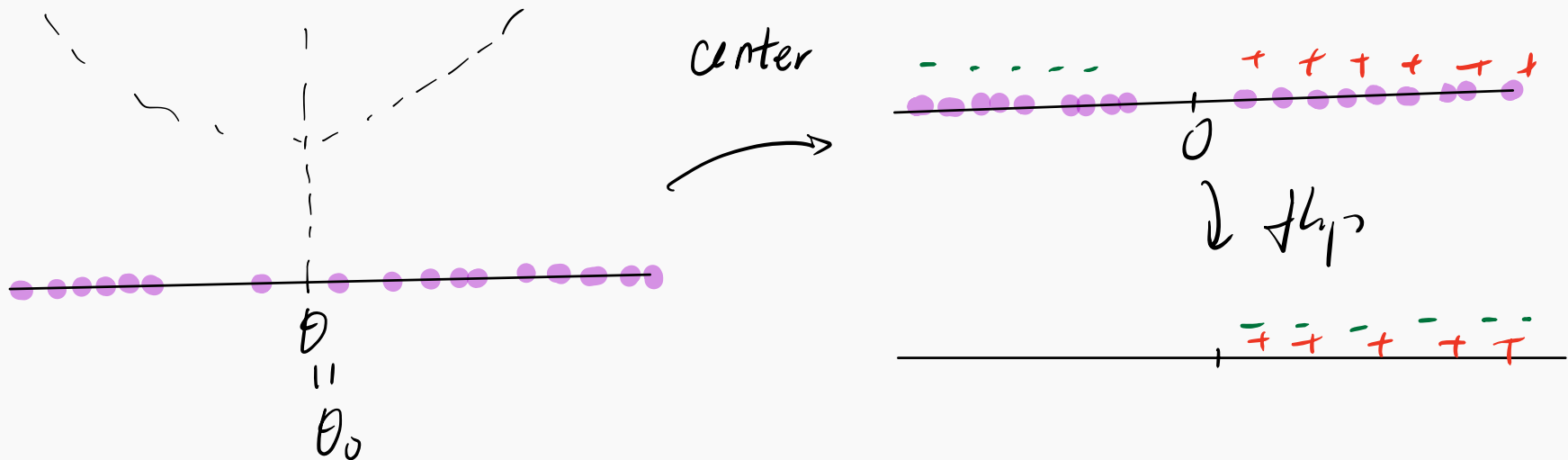
- Sign:

Define indicator variable for signs $\psi_i = \begin{cases} 1, & \text{if } Z_i > 0 \\ 0, & \text{if } Z_i < 0 \end{cases}$

Define positive signed rank of i th observation $\psi_i R_i = \begin{cases} R_i, & \text{if } Z_i > 0 \\ 0, & \text{if } Z_i < 0 \end{cases}$

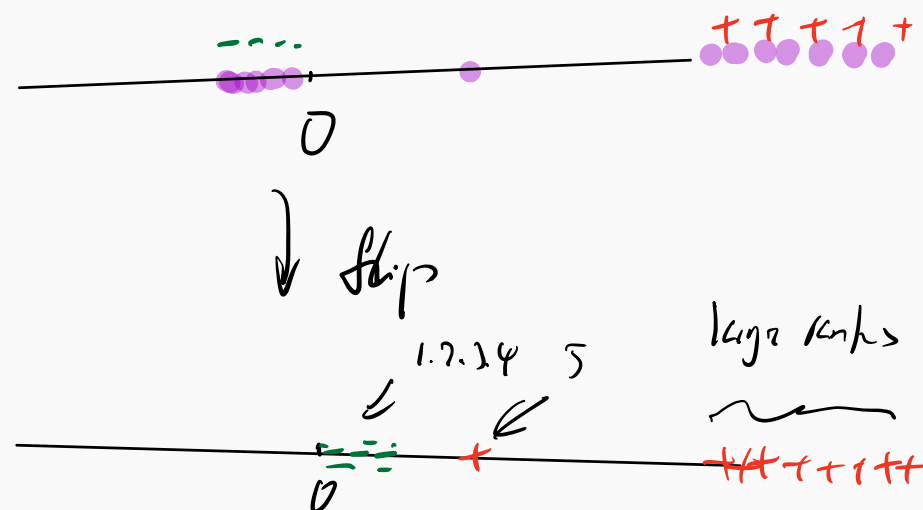
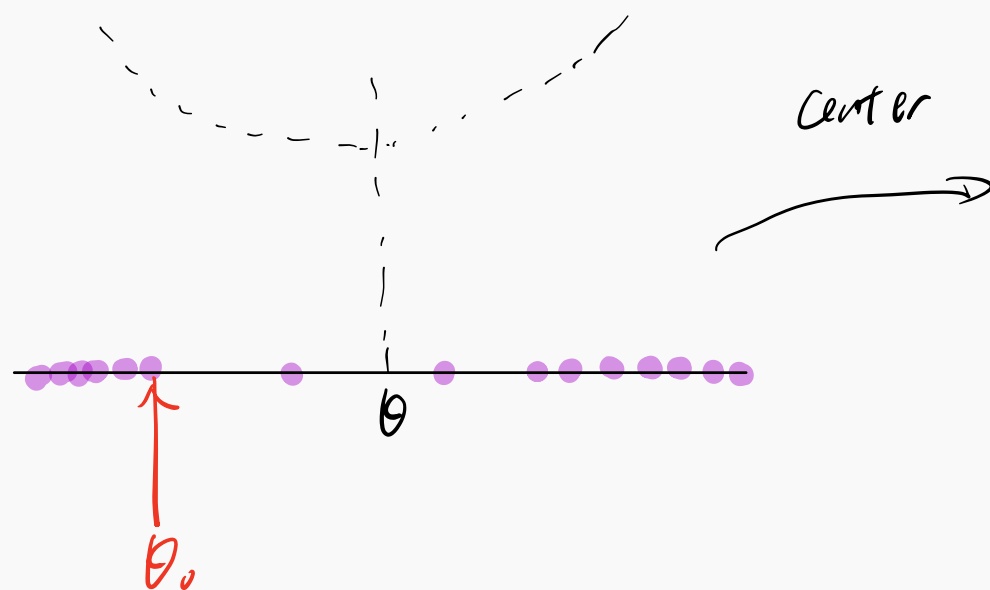
Motivation

- When the true unknown median θ is equal to the hypothesized value θ_0 :
 - the centered data will tend to be spread symmetrically around 0, and roughly half of the observations have positive signs, due to symmetric underlying distribution,
 - the ranks associated with negative and positive observations are roughly equal,
 - so the total ranks of those positive signed observation roughly is half of the total ranks of all observations.



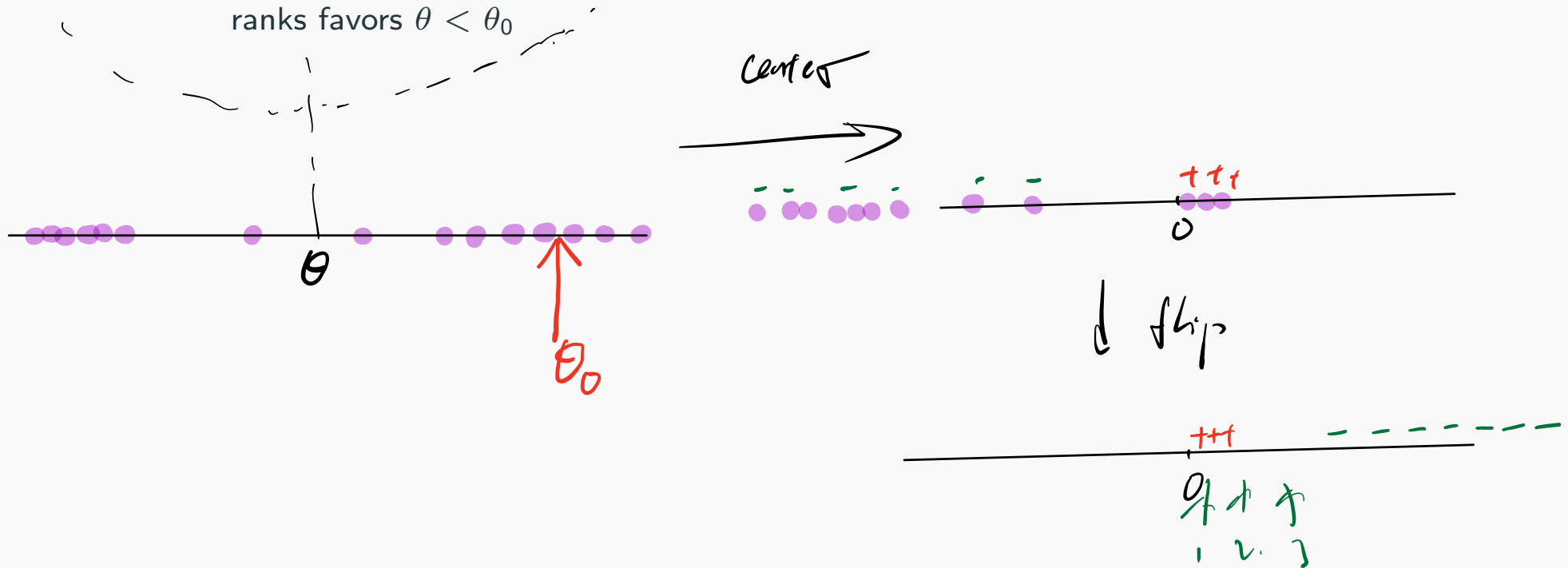
Motivation

- When the true unknown median $\theta > \theta_0$:
 - there will tend to be a larger portion of observations have positive signs and larger ranks associated with them,
 - so the total ranks of those positive signed observation is larger, so large total signed ranks favors $\theta > \theta_0$



Motivation

- When the true unknown median $\theta < \theta_0$:
 - there will tend to be a smaller portion of observations have positive signs and smaller ranks associated with them,
 - so the total ranks of those positive signed observation is smaller, so small total signed ranks favors $\theta < \theta_0$



Define Wilcoxon signed rank statistics

$$T^+ = \underbrace{\sum_{i=1}^n \psi_i R_i}_{\text{sum of positive signed ranks}}$$

Derivation of (exact) null distribution

When $H_0 : \theta = \theta_0$ is true :

- Think the process as randomly split ranks 1, 2, 3, ..., n into two groups
- T^+ is sum of one of groups

Because ranks and signs are independent under H_0

\Rightarrow each rank is equally likely to be + or - : $P(\psi_i = 1) = \frac{1}{2}$

\Rightarrow each configuration of permutation signed ranks $\{1\psi_1, 2\psi_2, \dots, n\psi_n\}$
occurs with $(\frac{1}{2})^n$

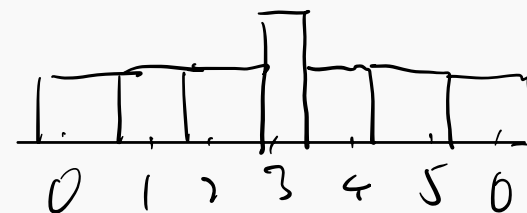
Example

$$n=3$$

$(1\psi_1, 2\psi_2, 3\psi_3)$	Prob. under H_0	T^+
$(0, 0, 0)$	$(\frac{1}{2})^3 = \frac{1}{8}$	0
$(1, 0, 0)$	$\frac{1}{8}$	1
$(0, 2, 0)$	$\frac{1}{8}$	2
$(0, 0, 3)$	$\frac{1}{8}$	3
$(1, 2, 0)$	$\frac{1}{8}$	3
$(1, 0, 3)$	$\frac{1}{8}$	4
$(0, 2, 3)$	$\frac{1}{8}$	5
$(1, 2, 3)$	$\frac{1}{8}$	6

\Rightarrow null dist. of T^+

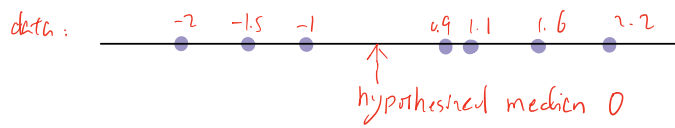
T^+	prob
0	$\frac{1}{8}$
1	$\frac{1}{8}$
2	$\frac{1}{8}$
3	$\frac{1}{8}$
4	$\frac{2}{8}$
5	$\frac{1}{8}$
6	$\frac{1}{8}$



We have derived the null distribution of T^+ without specifying the forms of the underlying populations beyond the point of requiring that they be continuous and symmetric about zero. This is why the test procedures based on T^+ are called **distribution-free procedures**.

From the null distribution of T^+ we can determine the critical value t_α and control the probability α of falsely rejecting H_0 when H_0 is true.

■ Wilcoxon signed rank test statistic

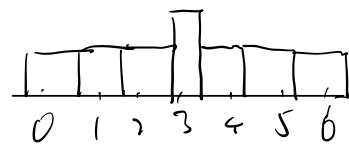
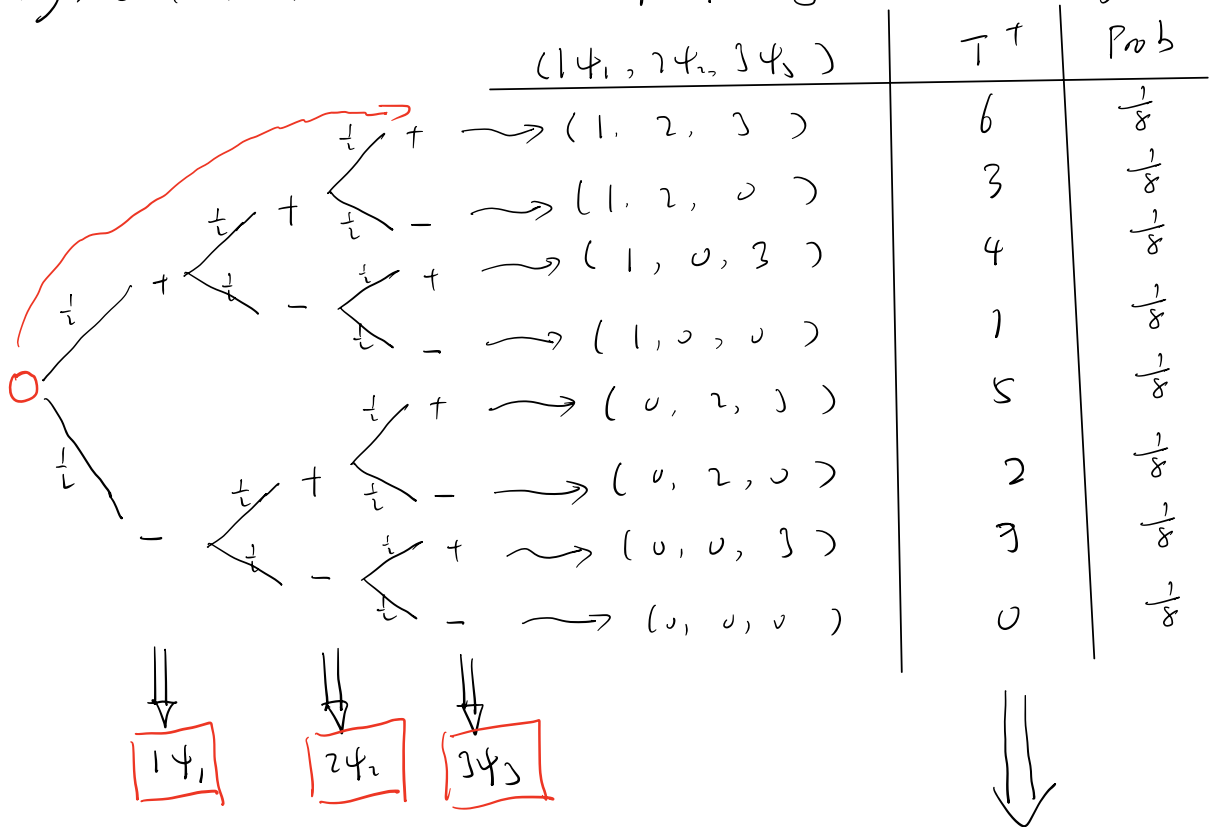


rank: 6 4 2 1 3 5 7

$$\Rightarrow T^+ = 1 + 3 + 5 + 7 = 16$$

■ (exer) null distribution of T^+

key: each rank 1, 2, 3, ..., n has equal probability to have + or - sign



null dist. of T^+

distribution-free property
= nonparametric

Large sample approximation of null distribution:

Reorder the data according to their absolute ranks, let $V_i = \psi_i S_i$ be the i th observation.

$$\begin{aligned} T^+ &= \sum_{i=1}^n \psi_i R_i \\ &= \sum_{i=1}^n V_i \end{aligned}$$

V_i are mutually independent dichotomous random variables with $P(V_i = 1) = P(V_i = 0) = \frac{1}{2}$. T^+ is sum of independent random variables, follows from standard theory for sums of mutually independent, but not identically distributed, random variables, such as the Liapounov central limit theorem (cf. Randles and Wolfe (1979, p. 423)), it has an asymptotic normality distribution.

Define standardized Wilcoxon signed rank statistics

$$\begin{aligned} \Rightarrow \text{standardized form } T^* &= \frac{T^+ - E_0(T^+)}{\{\text{var}_0(T^+)\}^{1/2}} = \frac{T^+ - \frac{n(n+1)}{4}}{\left\{ \frac{n(n+1)(2n+1)}{24} \right\}^{1/2}} \\ &\sim N(0, 1) \quad \text{if } n \text{ is large} \end{aligned}$$

n: sample size

Optima 1 :

$$\textcircled{D} \quad E_0(T^+) = E\left[\sum_{i=1}^n V_i\right] = \sum_{i=1}^n E[V_i]$$

$$E_0(V_i) = i\left(\frac{1}{2}\right) + 0\left(\frac{1}{2}\right) = \frac{i}{2}$$

$$\Rightarrow E_0(T^+) = \frac{1}{2} \sum_{i=1}^n i = \frac{1}{2} \left[\frac{n(n+1)}{2} \right] = \boxed{\frac{n(n+1)}{4}}$$

$$\textcircled{V} \quad \text{var}_0(T^+) = \text{var}\left(\sum_{i=1}^n V_i\right) = \sum_{i=1}^n \text{var}(V_i)$$

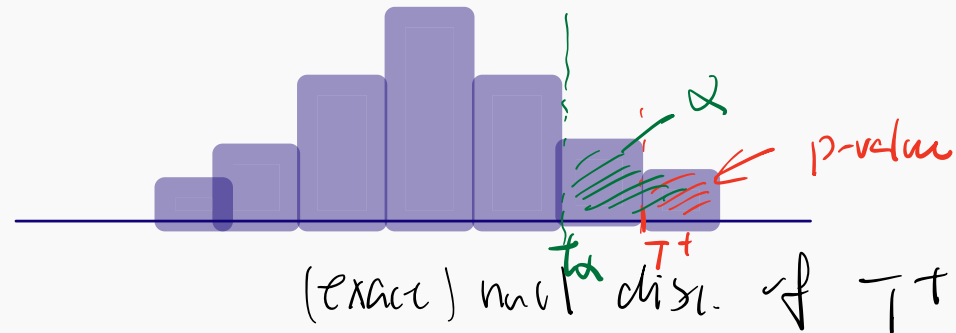
$$\text{var}_0(V_i) = E_0(V_i^2) - [E_0(V_i)]^2$$

$$= \left[i^2 \left(\frac{1}{2} \right) + 0^2 \left(\frac{1}{2} \right) \right] - \left[\frac{i}{2} \right]^2$$

$$= \frac{i^2}{2} - \frac{i^2}{4} = \frac{i^2}{4}$$

$$\Rightarrow \text{var}_0(T^+) = \frac{1}{4} \sum_{i=1}^n i^2 = \frac{1}{4} \left[\frac{n(n+1)(2n+1)}{6} \right] = \boxed{\frac{n(n+1)(2n+1)}{24}}$$

Procedure



a. One-Sided Upper-Tail Test.

To test

$$H_0 : \theta = 0$$

θ_0

versus

$$H_a : \theta > 0$$

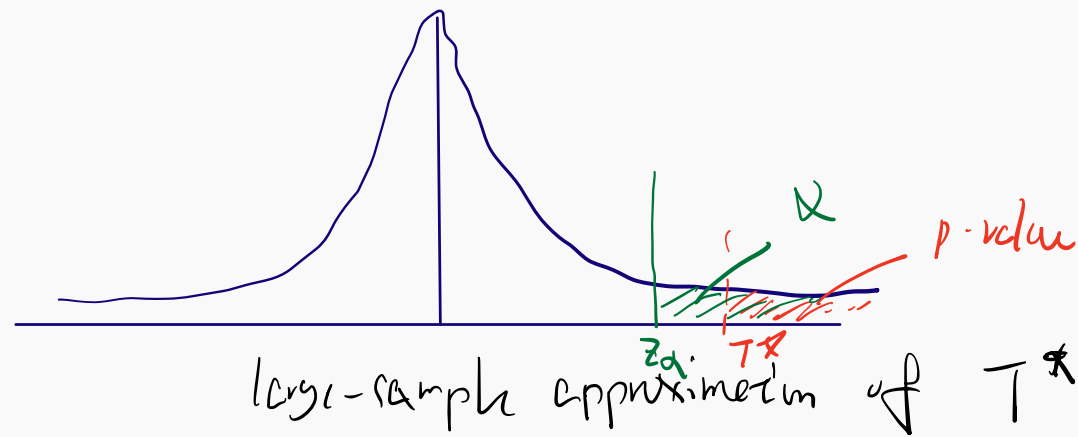
θ_0

Large T^+ / T^* favor H_a

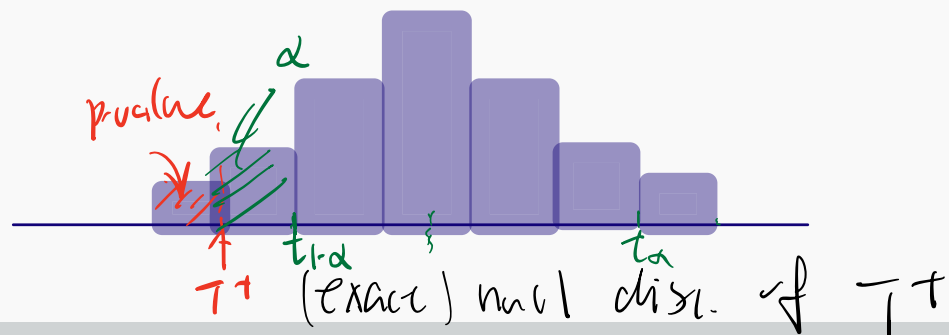
at the α level of significance, Reject H_0 if $T^+ \geq t_\alpha$; otherwise do not reject.

The normal approximation:

Reject H_0 if $T^* \geq z_\alpha$; otherwise do not reject.



Procedure



b. One-Sided Lower-Tail Test.

To test

$$H_0 : \theta = 0 \quad \text{vs} \quad \theta_0$$

versus

$$H_a : \theta < 0 \quad \text{vs} \quad \theta_0$$

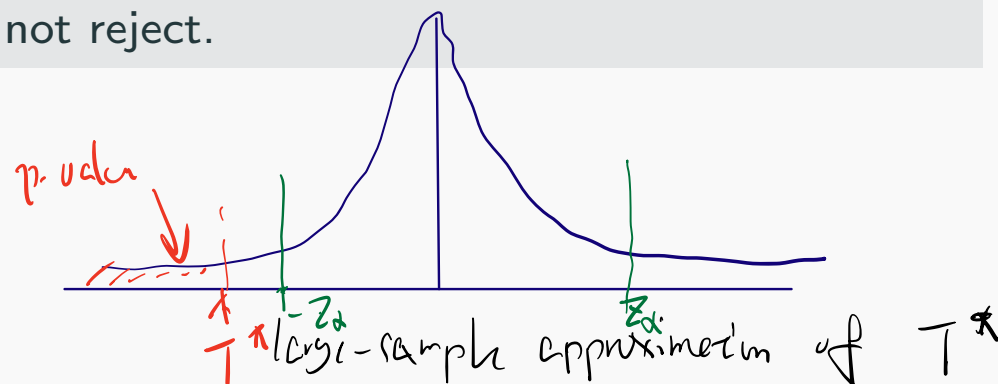
Small T^+ / T^- favor H_a

at the α level of significance, Reject H_0 if $T^+ \leq t_{1-\alpha} = \frac{n(n+1)}{2} - t_\alpha$; otherwise do not reject.

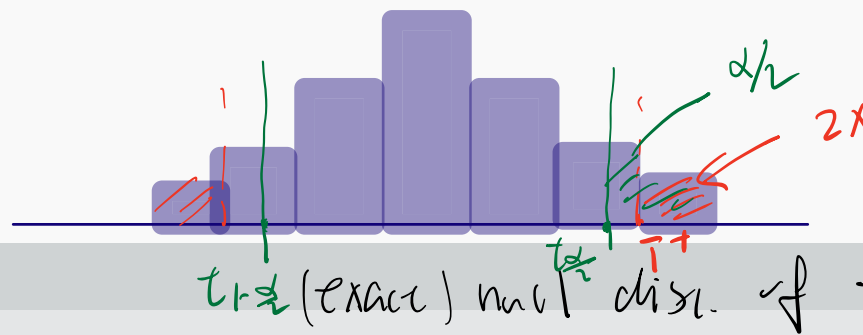
The normal approximation:

Reject H_0 if $T^* \leq -z_\alpha$; otherwise do not reject.

$$T^* = \frac{T^+ - \frac{n(n+1)}{2}}{\sqrt{\frac{n(n+1)}{12}}}$$



Procedure



c. Two-Sided Test.

To test

$$H_0 : \theta = 0$$

versus

$$H_a : \theta \neq 0$$

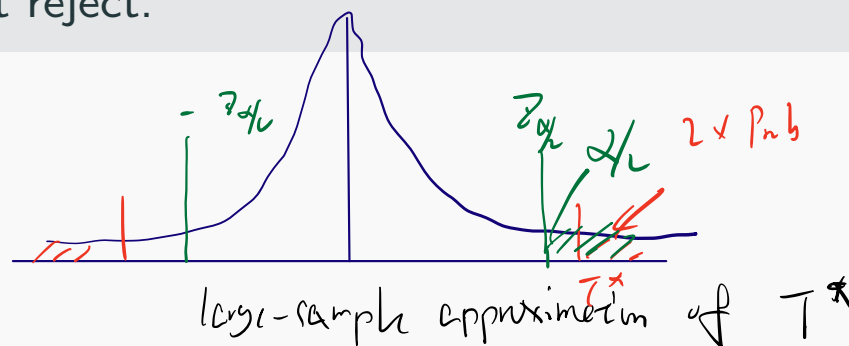
at the α level of significance, Reject H_0 if $T^+ \geq t_{\alpha/2}$ or $T^+ \leq t_{1-\alpha/2} = \frac{n(n+1)}{2} - t_{\alpha/2}$; otherwise do not reject.

Both small and large T^+/T^* favor H_a

This two-sided procedure is the two-sided symmetric test with $\alpha/2$ probability in each tail of the null distribution of T^+ .

The normal approximation:

Reject H_0 if $|T^*| \geq z_{\alpha/2}$; otherwise do not reject.

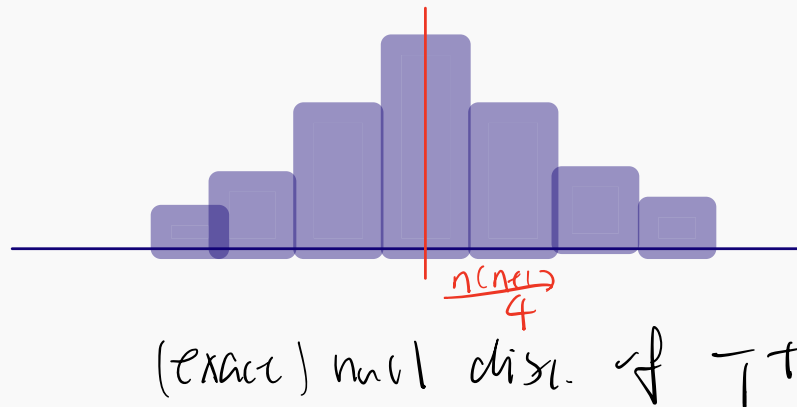


An estimator associated with Wilcoxon's signed rank statistics (Hodges-Lehmann)

The null distribution of the statistic T^+ is symmetric about its mean, $n(n+1)/4$. A natural estimator of θ is the amount that should be subtracted from each Z_i so that the value of T^+ , when applied to the shifted sample $X_1 - \hat{\theta}, \dots, X_n - \hat{\theta}$, is as close to $n(n+1)/4$ as possible.

Roughly speaking, we estimate θ by the amount ($\hat{\theta}$) that the Z sample should be shifted in order that $Z_1 - \hat{\theta}, \dots, Z_n - \hat{\theta}$ appears (when "viewed" by the signed rank statistic T^+) as a sample from a population with median 0.

If we have guessed correct median :



Optimal 1:

$$T^+ = \sum_{i=1}^n S_i R_i$$

$$S_i = I(X_i > \theta_0)$$

$$R_i = \sum_{j=1}^n I(|X_j - \theta_0| \leq |X_i - \theta_0|)$$

$$= \sum_{i=1}^n \sum_{j=1}^n I(|X_j - \theta_0| \leq |X_i - \theta_0|, X_i > \theta_0)$$

$$= \sum_{i=1}^n \sum_{j=1}^n I(|X_j - \theta_0| \leq X_i - \theta_0, X_i > \theta_0)$$

$$= \sum_{i=1}^n \sum_{j=1}^n I(\theta_0 - X_i \leq X_j - \theta_0 \leq X_i - \theta_0)$$

$$= \sum_{i=1}^n \sum_{j=1}^n I(2\theta_0 \leq X_j + X_i \leq 2X_i)$$

$$= \sum_{i=1}^n \sum_{j=1}^n I\left(\frac{X_i + X_j}{2} \geq \theta_0, X_j \leq X_i\right)$$

$$= \sum_{1 \leq i \leq j \leq n} 1 \left(\frac{X_i + X_j}{2} \geq \theta_0\right)$$

$$\Rightarrow \frac{1}{\frac{n(n+1)}{2}} \sum I\left(\frac{X_i + X_j}{2} \geq \theta_0\right) \approx \frac{1}{2}$$

	X_1	X_2	X_3
X_1	$\frac{X_1+X_1}{2}$	$\frac{X_1+X_2}{2}$	$\frac{X_1+X_3}{2}$
X_2		$\frac{X_2+X_2}{2}$	$\frac{X_2+X_3}{2}$
X_3			$\frac{X_3+X_3}{2}$

\Rightarrow The median θ_0 should be the value that splits the pairwise averages into 50% / 50%.

$$\approx \frac{n(n+1)}{4}$$

The Walsh Averages.

Each of the $n(n+1)/2$ averages $(X_i + X_j)/2, i \leq j = 1, \dots, n$, is called a Walsh average.

Estimate

To estimate the median θ , form the $M = n(n+1)/2$ averages $(X_i + X_j)/2$, for $i \leq j = 1, \dots, n$. The estimator of θ associated with the Wilcoxon signed rank statistic T^+ is

$$\hat{\theta} = \text{median} \left\{ \frac{\cancel{X_i} + \cancel{X_j}}{2}, i \leq j = 1, \dots, n \right\}.$$

Let $W^{(1)} \leq \dots \leq W^{(M)}$ denote the ordered values of $(X_i + X_j)/2$, where $M = \frac{n(n+1)}{2}$

- Then if M is odd, say $M = 2k + 1$, we have $k = (M - 1)/2$ and

$$\hat{\theta} = W^{(k+1)},$$

the value that occupies position $k + 1$ in the list of the ordered $(X_i + X_j)/2$ averages.

- If M is even, say $M = 2k$, then $k = M/2$ and

$$\hat{\theta} = \frac{W^{(k)} + W^{(k+1)}}{2}$$

That is, when M is even, $\hat{\theta}$ is the average of the two $(X_i + X_j)/2$ values that occupy positions k and $k + 1$ in the ordered list of the M $(X_i + X_j)/2$ averages. The $(X_i + X_j)/2$.

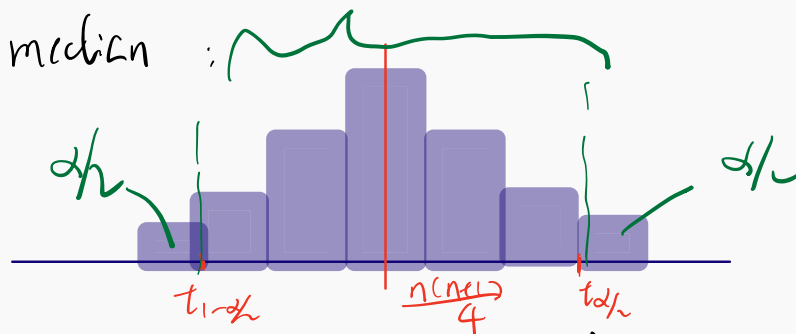
Confidence interval based on Wilcoxon's signed rank test

The true population median θ_0 is the value such that the number of Walsh averages above it is the Wilcoxon's signed rank statistics T^+ which should be centered at $\frac{n(n+1)}{4}$ with some natural variation. So we use the natural variation of T^+ reverse engineer the most probable region of θ_0 .

$$T^+ = \sum_{1 \leq i < j \leq n} I\left(\frac{X_i + X_j}{2} \geq \theta_0\right)$$

total number of Walsh averages above true median θ_0

If we have guessed correct median :



(exact) null dist. of T^+

$$P(t_{1-\alpha/2} \leq T^+ \leq t_{\alpha/2}) = 1-\alpha \approx 95\%$$

with probability $(1-\alpha)100\%$

\Rightarrow total number of Walsh averages above true median should be between $t_{1-\alpha/2}$ and $t_{\alpha/2}$

Procedure

$$W^{(1)} \quad \dots \quad W^{(\frac{n(n+1)}{2})}$$

$$\uparrow \quad W^{(\frac{n(n+1)}{2} + 1 - t_{\alpha/2})} \quad \uparrow \quad W^{(t_{\alpha/2})}$$

$(1 - \alpha)100\%$ Confidence Interval

For a symmetric two-sided confidence interval for median θ , with confidence coefficient $1 - \alpha$, set

$$\frac{t_{\alpha} - \frac{n(n+1)}{2} + 1 - t_{\alpha/2}}{2},$$

where $t_{\alpha/2}$ is the upper $(\alpha/2)$ th percentile point of the null distribution of T^+ .

The $100(1 - \alpha)\%$ confidence interval (θ_L, θ_U) for θ that is associated with the Wilcoxon signed rank statistics is

$$\theta_L = W^{(\frac{n(n+1)}{2} + 1 - t_{\alpha/2})}, \theta_U = W^{(M+1 - t_{\alpha/2})} = W^{(t_{\alpha/2})}$$

where $M = n(n+1)/2$ and $W^{(1)} \leq \dots \leq W^{(M)}$ are the ordered values of the $(X_i + X_j)/2$ Walsh averages. θ_L is the Walsh average that occupies position $t_{\alpha/2}$ in the list. The upper end point θ_U is the Walsh average that occupies the position $M+1 - t_{\alpha/2}$ in this ordered list.

$$W^{(1)} \quad W^{(2)} \quad W^{(3)} \quad W^{(4)} \quad W^{(5)} \quad W^{(6)} \quad W^{(7)} \quad W^{(8)}$$

$$t_{\alpha/2} = 7$$

$$\Rightarrow \text{C.I.} \therefore (\theta_L = W^{(2)}, \theta_U = W^{(7)})$$

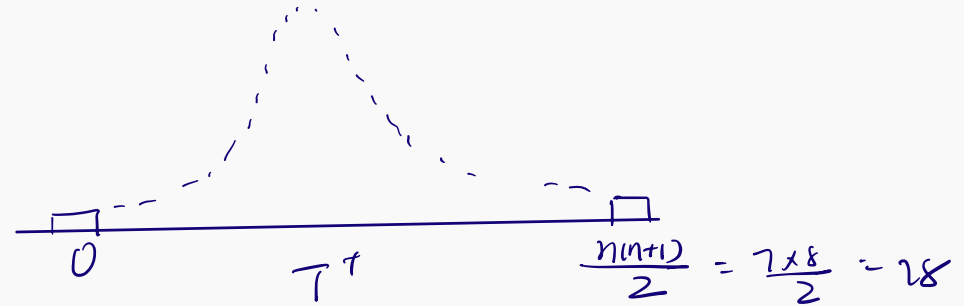
$$8+1-7 = 2$$

Example: The Mariner and the Pioneer Spacecraft Data

The data were reported by Anderson, Efron, and Wong (1970). The seven observations represent average measurements of, the ratio of the mass of the Earth to that of the moon, obtained from seven different spacecraft.

On the basis of the previous (2-3 years earlier) Ranger spacecraft findings, scientists had considered the value of the ratio of the mass of the Earth to that of the moon to be approximately 81.3035. Thus, we are interested in testing $H_0 : \theta = 81.3035$ versus the alternative $\theta \neq 81.3035$.

		Step 1	Step 2	Step 3
i	X_i	$X'_i = X_i - 81.3035$	$\text{Sign}(\psi_i)$	$\text{Rank}(R_i)$
1	81.3001	-.0034	-	6
2	81.3015	-.0020	-	2
3	81.3006	-.0029	-	4
4	81.3011	-.0024	-	3
5	81.2997	-.0038	-	7
6	81.3005	-.0030	-	5
7	81.3021	-.0014	-	1



Exact test:

$$T^+ = 0$$

$$p\text{-value} = 1/2^7 = 0.015625$$

$$P(T^+ \leq 0 \text{ or } T^+ \geq 28) = 2P(T^+ \leq 0)$$

$$= 2P(T^+ = 0)$$

$$= 2 \times \left(\frac{1}{2}\right)^7$$

Confirm with built-in function:

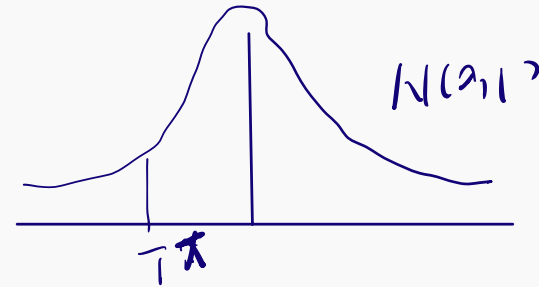
```
> wilcox.test(c(81.3001, 81.3015, 81.3006, 81.3011, 81.2997, 81.3005, 81.3021), mu=81.3035)
```

Wilcoxon signed rank exact test

data: c(81.3001, 81.3015, 81.3006, 81.3011, 81.2997, 81.3005, 81.3021)

V = 0, p-value = 0.01563

alternative hypothesis: true location is not equal to 81.3035













Large-sample approximation:

$$T^* = \frac{T^+ - \frac{n(n+1)}{4}}{\left\{ \frac{n(n+1)(2n+1)}{24} \right\}^{1/2}} = \frac{0 - [7(8)/4]}{[7(8)(15)/24]^{1/2}} = -2.366$$

$$p\text{-value} = 0.01798144 \Leftrightarrow 2P(Z \leq T^*) \quad \sim N(0,1) \quad > pnorm() \text{ in R}$$

Both the exact test and the large-sample approximation indicate the existence of strong evidence to reject the findings of the earlier Ranger spacecraft that $\theta = 81.3035$.

Walsh Energies	x_1	x_2	x_3	x_4
x_1				
x_2				
x_3				
x_4				

An estimate for median:

```
> library(Rfit)
> sort(walsh(c(81.3001,81.3015,81.3006,81.3011,81.2997,81.3005,81.3021)))
[1] 81.29970 81.29990 81.30010 81.30010 81.30015 81.30030 81.30035 81.30040
[9] 81.30050 81.30055 81.30060 81.30060 81.30060 81.30080 81.30080 81.30085
[17] 81.30090 81.30100 81.30105 81.30110 81.30110 81.30130 81.30130 81.30135
[25] 81.30150 81.30160 81.30180 81.30210
```

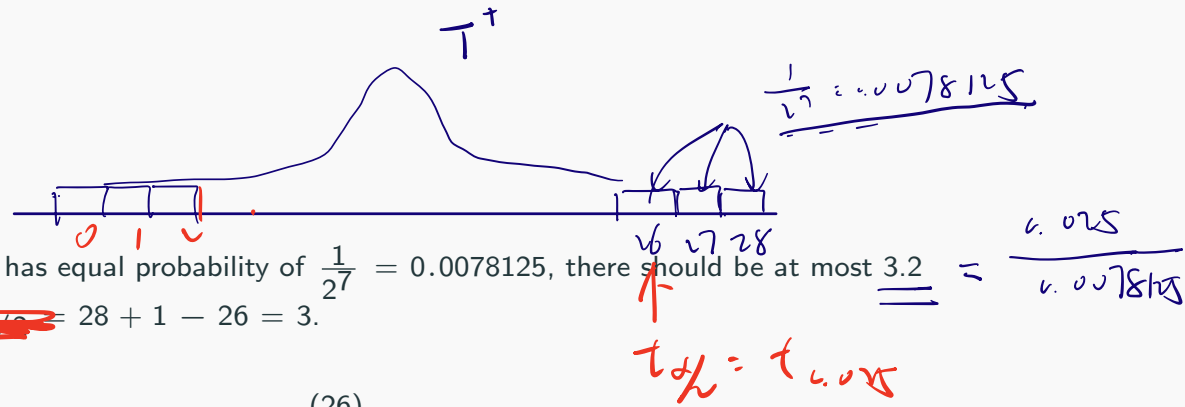
$M = 7(8)/2 = 28$, we see that $M = 2k$ with $k = 14$

$$\Rightarrow \hat{\theta} = \frac{w^{(14)} + w^{(15)}}{2} = \frac{81.3008 + 81.3008}{2} = 81.3008$$

95%

Confidence interval for median:

With $n = 7$ and $\alpha = .05$, each configuration under null has equal probability of $\frac{1}{2^7} = 0.0078125$, there should be at most $\underline{3.2} = \frac{0.025}{0.0078125}$ configurations to the right of $t_{\alpha/2} = 26$. Thus, ~~$t_{\alpha/2} = 28 + 1 - 26 = 3$~~



$$\theta_L = W^{(3)} = 81.3001 \text{ and } \theta_U = W^{(26)} = 81.3016$$

so that our 95% confidence interval for θ is

$$(\theta_L, \theta_U) = (81.3001, 81.3016)$$

Confirm with built-in function:

Correction \rightarrow
 ~~$m=28$~~
 ~~n~~

```
> wilcox.test(c(81.3001, 81.3015, 81.3006, 81.3011, 81.2997, 81.3005, 81.3021),  
m=28, n=81.3035, exact=T, conf.int=T, conf.level=0.95)
```

Wilcoxon signed rank exact test

```
data: c(81.3001, 81.3015, 81.3006, 81.3011, 81.2997, 81.3005, 81.3021)
```

```
V = 28, p-value = 0.01563
```

```
alternative hypothesis: true location is not equal to 28
```

```
95 percent confidence interval:
```

```
81.3001 81.3016
```

```
sample estimates:
```

```
(pseudo)median
```

```
81.3008
```

$W^{(1)}$ $W^{(2)}$ $W^{(3)}$

$W^{(26)}$ $W^{(27)}$ $W^{(28)}$