

Math 362: Mathematical Statistics II

Le Chen

le.chen@emory.edu

Emory University
Atlanta, GA

Last updated on April 28, 2021

2021 Spring

Chapter 14. Nonparametric Statistics

§ 14.1 Introduction

§ 14.2 The Sign Test

§ 14.3 Wilcoxon Tests

§ 14.4 The Kruskal-Wallis Test

§ 14.5 The Friedman Test

§ 14.6 Testing for Randomness

Chapter 14. Nonparametric Statistics

§ 14.1 Introduction

§ 14.2 The Sign Test

§ 14.3 Wilcoxon Tests

§ 14.4 The Kruskal-Wallis Test

§ 14.5 The Friedman Test

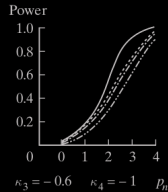
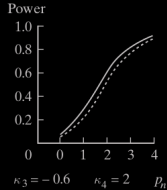
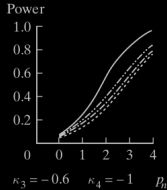
§ 14.6 Testing for Randomness

Nonparametric statistics

- ▶ *Distribution-free methods* : do not rely on assumptions that the data are drawn from a given parametric family of probability distributions.
- ▶ *Nonparametric statistics*: a statistic is defined to be a function on a sample and there is no dependency on any parameters, such as
 - *Order statistics*

Nonparametric vs. Parametric methods

– Power of Test



- Solid line: one-sample t-test (parametric test)
- Dashed lines: the sign test (nonparametric test)

Nonparametric vs. Parametric methods

Nonparametric methods usually produce

- ▶ Greater variance in point estimation
- ▶ Less power in hypothesis-testing
- ▶ Wider confidence intervals
- ▶ Lower probability of correct selection (in ranking and selection)
- ▶ Higher risk (in decision theory)

Hence, use nonparametric methods only when

The underlying assumptions for the
probability distributions are seriously doubtful.

Chapter 14. Nonparametric Statistics

§ 14.1 Introduction

§ 14.2 The Sign Test

§ 14.3 Wilcoxon Tests

§ 14.4 The Kruskal-Wallis Test

§ 14.5 The Friedman Test

§ 14.6 Testing for Randomness

- Let $\tilde{\mu}$ be the **median** of some unknown continuous pdf $f_Y(y)$:

$$\mathbb{P}(Y \leq \tilde{\mu}) = \mathbb{P}(Y \geq \tilde{\mu}) = \frac{1}{2}.$$

- For a random sample of size n is taken from $f_Y(y)$, in order to test

$$H_0 : \tilde{\mu} = \tilde{\mu}_0 \quad \text{vs} \quad H_0 : \tilde{\mu} \neq \tilde{\mu}_0,$$

let

$X :=$ the number of observations exceeding $\tilde{\mu}_0$

\Downarrow

1. $X \sim \text{Binomial}(n, 1/2)$.
2. Moreover, if n is large, by CLT,

$$\frac{X - \mathbb{E}[X]}{\sqrt{\text{Var}(X)}} = \frac{X - \frac{n}{2}}{\sqrt{n/4}} \underset{\text{aprox.}}{\sim} N(0, 1)$$

Sign test for median of a single sample

- ▶ When sample size n is large:

Let y_1, y_2, \dots, y_n be a random sample of size n from any continuous distribution having median $\tilde{\mu}$, where $n \geq 10$. Let k denote the number of y_i 's greater than $\tilde{\mu}_0$, and let $z = \frac{k-n/2}{\sqrt{n/4}}$.

- To test $H_0: \tilde{\mu} = \tilde{\mu}_0$ versus $H_1: \tilde{\mu} > \tilde{\mu}_0$ at the α level of significance, reject H_0 if $z \geq z_\alpha$.*
- To test $H_0: \tilde{\mu} = \tilde{\mu}_0$ versus $H_1: \tilde{\mu} < \tilde{\mu}_0$ at the α level of significance, reject H_0 if $z \leq -z_\alpha$.*
- To test $H_0: \tilde{\mu} = \tilde{\mu}_0$ versus $H_1: \tilde{\mu} \neq \tilde{\mu}_0$ at the α level of significance, reject H_0 if z is either (1) $\leq -z_{\alpha/2}$ or (2) $\geq z_{\alpha/2}$. □*

- ▶ When sample size n is small: use the exact distribution of binomial distribution.

E.g.1 In a healthy adults, the median pH for synovial fluid is 7.39.

A random sample of $n = 43$ is chosen and test

$$H_0 : \tilde{\mu} = 7.39 \quad \text{vs} \quad H_0 : \tilde{\mu} \neq 7.39, \quad \text{at } \alpha = 0.10.$$

Subject	Synovial Fluid pH	Subject	Synovial Fluid pH
HW	7.02	BG	7.34
AD	7.35	GL	7.22
TK	7.32	BP	7.32
EP	7.33	NK	7.40
AF	7.15	LL	6.99
LW	7.26	KC	7.10
LT	7.25	FA	7.30
DR	7.35	ML	7.21
VU	7.38	CK	7.33
SP	7.20	LW	7.28
MM	7.31	ES	7.35
DF	7.24	DD	7.24
LM	7.34	SL	7.36
AW	7.32	RM	7.09
BB	7.34	AL	7.32
TL	7.14	BV	6.95
PM	7.20	WR	7.35
JG	7.41	HT	7.36
DH	7.77	ND	6.60
ER	7.12	SJ	7.29
DP	7.45	BA	7.31
FF	7.28		

Sol 1. We first count how many samples exceeding the median (i.e., obtain the value of X)

Subject	Synovial Fluid pH	Subject	Synovial Fluid pH
HW	7.02	BG	7.34
AD	7.35	GL	7.22
TK	7.32	BP	7.32
EP	7.33	NK	7.40
AF	7.15	LL	6.99
LW	7.26	KC	7.10
LT	7.25	FA	7.30
DR	7.35	ML	7.21
VU	7.38	CK	7.33
SP	7.20	LW	7.28
MM	7.31	ES	7.35
DF	7.24	DD	7.24
LM	7.34	SL	7.36
AW	7.32	RM	7.09
BB	7.34	AL	7.32
TL	7.14	BV	6.95
PM	7.20	WR	7.35
JG	7.41	HT	7.36
DH	7.77	ND	6.60
ER	7.12	SJ	7.29
DP	7.45	BA	7.31
FF	7.28		

4

1

2

3

Hence, we have $k = 4$, $n = 43$, and since n is large, we use the z test:

$$z = \frac{4 - 43/2}{\sqrt{43/4}} = -5.34.$$

Since the critical regions (two-sided test here) are

$$\begin{aligned} &(-\infty, -z_{\alpha/2}) \cup (z_{\alpha/2}, \infty) \\ &\quad || \\ &(-\infty, -2.58) \cup (2.58, \infty), \end{aligned}$$

we reject the hypothesis.

Or equivalently, the p -value is

$$2 \times \mathbb{P}(Z < -5.34) = 9.294658 \times 10^{-8}.$$



```
1 | > pnorm(-5.34) *2
2 | [1] 9.294658e-08
```

Sol 2. We can also carry out the exact computation thanks to computer:

The exact p-value should be

$$2 \times \mathbb{P}(X \leq 5) = 2 \sum_{k=0}^5 \binom{43}{k} \left(\frac{1}{2}\right)^{43} = 2.49951 \times 10^{-7},$$

which is smaller than $\alpha = 0.10$.

Hence, rejection!



```
1 | > pbinom(5,43,0.5) * 2
2 | [1] 2.49951e-07
```

Sign test for paired data

E.g. A manufacturer produces two products, A and B. The manufacturer wishes to know if consumers prefer product B over product A.

A sample of 10 consumers are each given product A and product B, and asked which product they prefer:

Preferences	Number
B	8
A	1
No preference	1

Test at $\alpha = 0.10$ that

H_0 : consumers do not prefer B over A

vs.

H_1 : consumers do prefer B over A.

Sol. We first remove the ties. So that we have a random (paired-data) sample of size $n = 9$.

Under H_0 , the consumers have no preference for B over A. Hence, we may believe that consumers will choose A or B with probability $\frac{1}{2}$.

Hence, to get more extreme values in this setting would give the p-value:

$$\mathbb{P}(X \geq 8) = \sum_{k=8}^9 \binom{9}{k} \left(\frac{1}{2}\right)^9 = 0.0195.$$

Conclusion, Rejection!



Chapter 14. Nonparametric Statistics

§ 14.1 Introduction

§ 14.2 The Sign Test

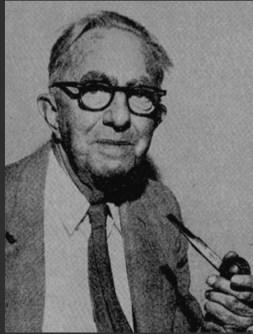
§ 14.3 Wilcoxon Tests

§ 14.4 The Kruskal-Wallis Test

§ 14.5 The Friedman Test

§ 14.6 Testing for Randomness

Frank Wilcoxon



Born	2 September 1892 County Cork, Ireland
Died	18 November 1965 (aged 73) Tallahassee, Florida, USA
Nationality	Irish American
Alma mater	Cornell University Rutgers University
Scientific career	
Fields	Chemistry Statistics
Institutions	American Cyanamid Company

Testing $H_0 : \mu = \mu_0$

Setup Let Y_1, \dots, Y_n be a set of independent variables with pdfs $f_{Y_1}(y), \dots, f_{Y_n}(y)$, respectively.

Assume that $f_{Y_i}(y)$ are continuous and symmetric.

Assume that all mean/median of f_{Y_i} are equal, denoted by μ .

Test $H_0 : \mu = \mu_0$ vs. $H_1 : \mu \neq \mu_0$.

Wilcoxon signed rank static

$$W = \sum_{k=1}^n R_k \mathbb{I}_{\{Y_k > \mu_0\}}$$

where R_i denotes the rank (increasing and starting from 1) of

$$\{|Y_1 - \mu_0|, |Y_2 - \mu_0|, \dots, |Y_n - \mu_0|\}$$

n	1	2	3
y_n	4.2	6.1	2.0
$y_n - 3.0$	1.2	3.1	-1.0
$ y_n - 3.0 $	1.2	3.1	1.0
r_n	2	3	1
$\mathbb{I}_{\{y_n > 3.0\}}$	1	1	0
$r_n \mathbb{I}_{\{y_n > 3.0\}}$	$u_2 = 2$	$u_3 = 3$	$u_1 = 0$

\Downarrow

$$w = 2 \times 1 + 3 \times 1 + 1 \times 0 = 5.$$

Let $\{y_1, \dots, y_n\}$ be For a sample of size n .

Some observations:

- ▶ r_i takes values in $\{1, 2, \dots, n\}$.
- ▶ w_i takes values in $\left\{0, 1, 2, \dots, \frac{n(n+1)}{2}\right\}$ with $1 + 2 + \dots + n = \frac{n(n+1)}{2}$.
- ▶ W is a discrete random variable:

w	0	1	\dots	$\frac{n(n+1)}{2}$
$\mathbb{P}(W = w)$				

Theorem Under the above setup and under H_0 ,

$$p_W(w) = \mathbb{P}(W = w) = \frac{c(w)}{2^n},$$

where $c(w)$ is the coefficient of e^{wt} in the expansion of

$$\prod_{k=1}^n (1 + e^{kt}).$$

Proof Under H_0 , $W = \sum_{k=1}^n U_k$ with follow the following distribution

$$U_k = \begin{cases} 0 & \text{with probability } 1/2 \\ k & \text{with probability } 1/2. \end{cases}$$

Then

$$M_W(t) = \prod_{k=1}^n M_{U_k}(t) = \prod_{k=1}^n \mathbb{E} \left(e^{U_k t} \right) = \prod_{k=1}^n \left(\frac{1}{2} + \frac{1}{2} e^{kt} \right).$$

Hence, we have


$$M_W(t) = \frac{1}{2^n} \prod_{k=1}^n \left(1 + e^{kt}\right).$$

On the other hand,

$$M_W(t) = \mathbb{E} \left(e^{Wt} \right) = \sum_{w=0}^{\frac{n(n+1)}{2}} e^{wt} p_W(w)$$

Equating the above two expressions, namely,

$$\frac{1}{2^n} \prod_{k=1}^n \left(1 + e^{kt}\right) = \sum_{w=0}^{\frac{n(n+1)}{2}} e^{wt} p_W(w),$$

proves the theorem. 

E.g. Find the pdf of W when $n = 2$ and 4.

Sol. When $n = 2$,

$$\begin{aligned}M_W(t) &= \frac{1}{2^2} (1 + e^t) (1 + e^{2t}) \\&= \frac{1}{2^2} (1 + e^t + e^{2t} + e^{3t}).\end{aligned}$$

Hence,

w	0	1	2	3
$p_W(w)$	1/4	1/4	1/4	1/4

When $n = 4$,

$$\begin{aligned} M_W(t) &= \frac{1}{2^4} (1 + e^t) (1 + e^{2t}) (1 + e^{3t}) (1 + e^{4t}) \\ &= \frac{1}{16} (e^{10t} + e^{9t} + e^{8t} + 2e^{7t} + 2e^{6t} + 2e^{5t} + 2e^{4t} + 2e^{3t} + e^{2t} + e^t + 1) \end{aligned}$$

Hence,

w	0	1	2	3	4	5	6	7	8	9	10
$p_W(w)$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{2}{16}$	$\frac{2}{16}$	$\frac{2}{16}$	$\frac{2}{16}$	$\frac{2}{16}$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$



```

1 sage: var('k,t')
2 (k, t)
3 sage: product(1+e^(k*t),k,1,4)
4 e^(10*t) + e^(9*t) + e^(8*t) + 2*e^(7*t) + 2*e^(6*t) + 2*e^(5*t) + 2*e^(4*t) + 2*
  e^(3*t) + e^(2*t) + e^t + 1

```


E.g. Shark studies:

Table 14.3.2 Measurements Made on Ten Sharks Caught Near Santa Catalina		
Total Length (mm)	Height of First Dorsal Fin (mm)	TL/HDI
906	68	13.32
875	67	13.06
771	55	14.02
700	59	11.86
869	64	13.58
895	65	13.77
662	49	13.51
750	52	14.42
794	55	14.44
787	51	15.43

Past data show that the true average TL/HDI ratio should be 14.60.

Let $Y_i = TL/HDI$.

Does the data support the above claim, namely, test

$$H_0 : \mu = 14.60 \quad \text{vs.} \quad H_1 : \mu \neq 14.60.$$

Set $\alpha = 0.05$.

Sol. Computing the Wilcoxon signed rank statistics:

Table 14.3.3 Computations for Wilcoxon Signed Rank Test					
$TL/HDI (= y_i)$	$y_i - 14.60$	$ y_i - 14.60 $	r_i	z_i	$r_i z_i$
13.32	-1.28	1.28	8	0	0
13.06	-1.54	1.54	9	0	0
14.02	-0.58	0.58	3	0	0
11.86	-2.74	2.74	10	0	0
13.58	-1.02	1.02	6	0	0
13.77	-0.83	0.83	4.5	0	0
13.51	-1.09	1.09	7	0	0
14.42	-0.18	0.18	2	0	0
14.44	-0.16	0.16	1	0	0
15.43	+0.83	0.83	4.5	1	4.5

Hence, $w = 4.5$.

Now check the table to find the critical region:

$$C = \{w : w \leq 8 \quad \text{or} \quad w \geq 47\}.$$

Conclusion: Rejection!



```
1 > x <- c(13.32, 13.06, 14.02, 11.86, 13.58, 13.77, 13.51, 14.42, 14.44, 15.43)
2 > wilcox.test(x, mu = 14.60, alternative = "two.sided")
3
4     Wilcoxon signed rank exact test
5
6 data: x
7 V = 15, p-value = 0.123
8 alternative hypothesis: true location is not equal to 14.6
```

Large-sample Wilcoxon Signed Rank Test

Theorem Under the same setup and H_0 , we have

$$\mathbb{E}(W) = \frac{n(n+1)}{4} \quad \text{and} \quad \text{Var}(W) = \frac{n(n+1)(2n+1)}{24}.$$

Proof.

$$\begin{aligned} \mathbb{E}(W) &= \mathbb{E}\left(\sum_{k=1}^n U_k\right) = \sum_{k=1}^n \left(0 \cdot \frac{1}{2} + k \cdot \frac{1}{2}\right) \\ &= \sum_{k=1}^n \frac{k}{2} = \frac{n(n+1)}{4}. \end{aligned}$$

$$\begin{aligned} \text{Var}(W) &= \text{Var}\left(\sum_{k=1}^n U_k\right) = \sum_{k=1}^n \text{Var}(U_k) = \sum_{k=1}^n [\mathbb{E}(U_k^2) - \mathbb{E}(U_k)^2] \\ &= \sum_{k=1}^n \left[\frac{k^2}{2} - \left(\frac{k}{2}\right)^2\right] = \sum_{k=1}^n \frac{k^2}{4} = \frac{1}{4} \frac{n(n+1)(2n+1)}{6} \end{aligned}$$

Hence when n is large (usually $n \geq 12$),

$$\frac{W - \mathbb{E}(W)}{\sqrt{\text{Var}(W)}} = \frac{W - [n(n+1)]/4}{\sqrt{[n(n+1)(2n+1)]/24}} \stackrel{\text{approx}}{\sim} N(0, 1).$$

\Downarrow

Let w be the signed rank statistic based on n independent observations, each drawn from a continuous and symmetric pdf, where $n > 12$. Let

$$z = \frac{w - [n(n+1)]/4}{\sqrt{[n(n+1)(2n+1)]/24}}$$

- a. To test $H_0: \mu = \mu_0$ versus $H_1: \mu > \mu_0$ at the α level of significance, reject H_0 if $z \geq z_\alpha$.*
- b. To test $H_0: \mu = \mu_0$ versus $H_1: \mu < \mu_0$ at the α level of significance, reject H_0 if $z \leq -z_\alpha$.*
- c. To test $H_0: \mu = \mu_0$ versus $H_1: \mu \neq \mu_0$ at the α level of significance, reject H_0 if z is either $(1) \leq -z_{\alpha/2}$ or $(2) \geq z_{\alpha/2}$. □*

The Wilcoxon Rank Sum Test

– Nonparametric counterpart of the pooled two-sample t-test

Setup Let x_1, \dots, x_n and y_{n+1}, \dots, y_{n+m} be two independent random samples from $f_X(x)$ and $f_Y(y)$, respectively.

Assume that $f_X(x)$ and $f_Y(y)$ are the same except for a possible shift in location.

Test $H_0 : \mu_X = \mu_Y$ vs. ...

Test statistic

$$W = \sum_{k=1}^{n+m} R_k Z_k$$

where R_i is the rank (starting from the lowest with rank 1) and

$$Z_i = \begin{cases} 1 & \text{the } i\text{th entry comes from } f_X(x) \\ 0 & \text{the } i\text{th entry comes from } f_Y(y). \end{cases}$$

Theorem Under the above setup and under H_0 ,

$$\mathbb{E}[W] = \frac{n(n+m+1)}{2} \quad \text{and} \quad \text{Var}(W) = \frac{nm(n+m+1)}{12}.$$

Hence when sample sizes are large, namely, $n, m > 10$,

$$\frac{W - \mathbb{E}(W)}{\sqrt{\text{Var}(W)}} = \frac{W - [n(n+m+1)]/2}{\sqrt{[nm(n+m+1)]/12}} \underset{\sim}{\text{approx}} N(0, 1).$$

E.g. Baseball ...

Test if $H_0 : \mu_X = \mu_Y$ vs. $H_0 : \mu_X \neq \mu_Y$

Obs. #	Team	Time (min)	r_i	z_i	$r_i z_i$
1	Baltimore	177	21	1	21
2	Boston	177	21	1	21
3	California	165	7.5	1	7.5
4	Chicago (AL)	172	14.5	1	14.5
5	Cleveland	172	14.5	1	14.5
6	Detroit	179	24.5	1	24.5
7	Kansas City	163	5	1	5
8	Milwaukee	175	18	1	18
9	Minnesota	166	9.5	1	9.5
10	New York (AL)	182	26	1	26
11	Oakland	177	21	1	21
12	Seattle	168	12.5	1	12.5
13	Texas	179	24.5	1	24.5
14	Toronto	177	21	1	21
15	Atlanta	166	9.5	0	0
16	Chicago (NL)	154	1	0	0
17	Cincinnati	159	2	0	0
18	Houston	168	12.5	0	0
19	Los Angeles	174	16.5	0	0
20	Montreal	174	16.5	0	0
21	New York (NL)	177	21	0	0
22	Philadelphia	167	11	0	0
23	Pittsburgh	165	7.5	0	0
24	San Diego	161	3.5	0	0
25	San Francisco	164	6	0	0
26	St. Louis	161	3.5	0	0

Group X

Group Y

$w' = 240.5$

In this case, $n = 14$, $m = 12$, $w = 240.5$.

$$\mathbb{E}(W) = \frac{14(14 + 12 + 1)}{2} = 189,$$

$$\text{Var}(W) = \frac{14 \times 12 \times (14 + 12 + 1)}{12} = 378.$$

Hence, the approximate z-score is

$$z = \frac{w - \mathbb{E}(W)}{\sqrt{\text{Var}(W)}} = \frac{240.5 - 189}{\sqrt{378}} = 2.65.$$

...



Chapter 14. Nonparametric Statistics

§ 14.1 Introduction

§ 14.2 The Sign Test

§ 14.3 Wilcoxon Tests

§ 14.4 The Kruskal-Wallis Test

§ 14.5 The Friedman Test

§ 14.6 Testing for Randomness

The Kruskal-Wallis Test

What is the nonparametric counterpart for the one-way ANOVA?

Setup Suppose that $k \geq 2$ independent sample of size n_1, \dots, n_k are drawn from k

identically shaped and scaled pdfs,
except for possibly different medians.

Let $\tilde{\mu}_1, \dots, \tilde{\mu}_k$ be the medians.

Test $H_0 : \tilde{\mu}_1 = \tilde{\mu}_2 = \dots = \tilde{\mu}_k$ vs. $H_1 : \text{not all the } \tilde{\mu}_i\text{'s are equal.}$

Remark This is the test for median not mean, but if pdfs are symmetric, they are the same.

Kruskal-Wallis statistic B

$$B = \frac{12}{n(n+1)} \sum_{j=1}^k \frac{R_{.j}^2}{n_j} - 3(n+1)$$

where

Table 14.4.1 Notation for Kruskal-Wallis Procedure				
	<i>Treatment Level</i>			
	1	2	...	k
	$Y_{11}(R_{11})$	$Y_{12}(R_{12})$		$Y_{1k}(R_{1k})$
	$Y_{21}(R_{21})$			
	\vdots	\vdots	\dots	\vdots
	$Y_{n_1 1}(R_{n_1 1})$	$Y_{n_2 2}(R_{n_2 2})$		$Y_{n_k k}(R_{n_k k})$
Totals	$R_{.1}$	$R_{.2}$		$R_{.k}$

Theorem Under the above setup and under H_0 , then

$$B = \frac{12}{n(n+1)} \sum_{j=1}^k \frac{R_j^2}{n_j} - 3(n+1) \underset{\sim}{\approx} \chi_{k-1}^2.$$

H_0 should be rejected at the α level of significance if $b > \chi_{1-\alpha, k-1}^2$.

E.g. Lottery over the year 1969; Whether lottery is random?

Test if $H_0 : \tilde{\mu}_{\text{Jan}} = \tilde{\mu}_{\text{Feb}} = \cdots = \tilde{\mu}_{\text{Dec}}$ at $\alpha = 0.01$


Table 14.4.2 1969 Draft Lottery, Highest Priority (001) to Lowest Priority (366)												
Date	Jan.	Feb.	Mar.	Apr.	May	June	July	Aug.	Sept.	Oct.	Nov.	Dec.
1	305	086	108	032	330	249	093	111	225	359	019	129
2	159	144	029	271	298	228	350	045	161	125	034	328
3	251	297	267	083	040	301	115	261	049	244	348	157
4	215	210	275	081	276	020	279	145	232	202	266	165
5	101	214	293	269	364	028	188	054	082	024	310	056
6	224	347	139	253	155	110	327	114	006	087	076	010
7	306	091	122	147	035	085	050	168	008	234	051	012
8	199	181	213	312	321	366	013	048	184	283	097	105
9	194	338	317	219	197	335	277	106	263	342	080	043
10	325	216	323	218	065	206	284	021	071	220	282	041
11	329	150	136	014	037	134	248	324	158	237	046	039
12	221	068	300	346	133	272	015	142	242	072	066	314
13	318	152	259	124	295	069	042	307	175	138	126	163
14	238	004	354	231	178	356	331	198	001	294	127	026
15	017	089	169	273	130	180	322	102	113	171	131	320
16	121	212	166	148	055	274	120	044	207	254	107	096
17	235	189	033	260	112	073	098	154	255	288	143	304
18	140	292	332	090	278	341	190	141	246	005	146	128
19	058	025	200	336	075	104	227	311	177	241	203	240
20	280	302	239	345	183	360	187	344	063	192	185	135
21	186	363	334	062	250	060	027	291	204	243	156	070
22	337	290	265	316	326	247	153	339	160	117	009	053
23	118	057	256	252	319	109	172	116	119	201	182	162
24	059	236	258	002	031	358	023	036	195	196	230	095
25	052	179	343	351	361	137	067	286	149	176	132	084
26	092	365	170	340	357	022	303	245	018	007	309	173
27	355	205	268	074	296	064	289	352	233	264	047	078
28	077	299	223	262	308	222	088	167	257	094	281	123
29	349	285	362	191	226	353	270	061	151	229	099	016
30	164		217	208	103	209	287	333	315	038	174	003
31	211		030		313		193	011		079		100
Totals:	6236	5886	7000	6110	6447	5872	5628	5377	4719	5656	4462	3768

Sol. Rank the lottery for the year (see the previous table).

Compute b using the formula:

$$\begin{aligned} b &= \frac{12}{366 \times 367} \left[\frac{6236^2}{31} + \frac{5886^2}{29} + \cdots + \frac{3768^2}{31} \right] - 3 \times 367 \\ &= 25.95. \end{aligned}$$

Critical region is $C = \{b : b \geq \chi_{0.99,11}^2 = 24.725\}$.

Conclusion: Reject (Lottery is *NOT* random). 

Chapter 14. Nonparametric Statistics

§ 14.1 Introduction

§ 14.2 The Sign Test

§ 14.3 Wilcoxon Tests

§ 14.4 The Kruskal-Wallis Test

§ 14.5 The Friedman Test

§ 14.6 Testing for Randomness

The Friedman Test

What is the nonparametric counterpart for the two-way ANOVA?

Setup Suppose that $k \geq 2$ independent sample of size n_1, \dots, n_k are drawn from k

identically shaped and scaled pdfs,
except for possibly different medians.

Assume that $n_1 = \dots = n_k$.

Samples can be further partitioned into b blocks.

Let $\tilde{\mu}_1, \dots, \tilde{\mu}_k$ be the medians.

Test $H_0 : \tilde{\mu}_1 = \tilde{\mu}_2 = \dots = \tilde{\mu}_k$ vs. $H_1 : \text{not all the } \tilde{\mu}_i\text{'s are equal.}$

Remark This is the test for median not mean, but if pdfs are symmetric, they are the same.

The Friedman Test Statistic:

Reject H_0 at the α level if

$$G = \frac{12}{bk(k+1)} \sum_{j=1}^k R_{.j}^2 - 3b(k+1) \geq \chi_{1-\alpha, k-1}^2.$$

where $R_{.j}$ is the within-block ranks.

E.g. Baseball ...

Test if $H_0 : \tilde{\mu}_{\text{Narrow}} = \tilde{\mu}_{\text{Wide}}$ at $\alpha = 0.01$

Table 14.5.1 Times (sec) Required to Round First Base				
Player	Narrow-Angle	Rank	Wide-Angle	Rank
1	5.50	1	5.55	2
2	5.70	1	5.75	2
3	5.60	2	5.50	1
4	5.50	2	5.40	1
5	5.85	2	5.70	1
6	5.55	1	5.60	2
7	5.40	2	5.35	1
8	5.50	2	5.35	1
9	5.15	2	5.00	1
10	5.80	2	5.70	1
11	5.20	2	5.10	1
12	5.55	2	5.45	1
13	5.35	1	5.45	2
14	5.00	2	4.95	1
15	5.50	2	5.40	1
16	5.55	2	5.50	1
17	5.55	2	5.35	1
18	5.50	1	5.55	2
19	5.45	2	5.25	1
20	5.60	2	5.40	1
21	5.65	2	5.55	1
22	6.30	2	6.25	1
		39		27

Sol. $k = 2$, $b = 22$

Compute the rank within each block (see the previous table)

Compute the g statistic:


$$g = \frac{12}{22 \times 2 \times (2 + 1)} [39^2 + 27^2] - 3 \times 22 \times (2 + 1) = \frac{72}{11} \approx 6.54.$$

Critical region is

$$C = \{g : g \geq \chi_{0.95,1}^2 = 3.84\}.$$

The p -value is

$$\mathbb{P}\left(\chi_1^2 \geq \frac{72}{11}\right) = 0.01051525.$$

Conclusion: Reject. 

R Code for this problem:

```
1 C1 <- c(
2 5.50, 5.70, 5.60, 5.50, 5.85, 5.55, 5.40, 5.50, 5.15, 5.80, 5.20,
3 5.55, 5.35, 5.00, 5.50, 5.55, 5.55, 5.50, 5.45, 5.60, 5.65, 6.30)
4 C2 <- c(
5 5.55, 5.75, 5.50, 5.40, 5.70, 5.60, 5.35, 5.35, 5.00, 5.70, 5.10,
6 5.45, 5.45, 4.95, 5.40, 5.50, 5.35, 5.55, 5.25, 5.40, 5.55, 6.25)
7 angles <- matrix(
8   cbind(C1, C2),
9   nrow = 22,
10  byrow = FALSE,
11  dimnames = list(1:22, c("Narrow", "Wide")))
12 )
13 friedman.test(angles)
```

Here is the output:

```
1 > C1 <- c(
2 + 5.50, 5.70, 5.60, 5.50, 5.85, 5.55, 5.40, 5.50, 5.15, 5.80, 5.20,
3 + 5.55, 5.35, 5.00, 5.50, 5.55, 5.55, 5.50, 5.45, 5.60, 5.65, 6.30)
4 > C2 <- c(
5 + 5.55, 5.75, 5.50, 5.40, 5.70, 5.60, 5.35, 5.35, 5.00, 5.70, 5.10,
6 + 5.45, 5.45, 4.95, 5.40, 5.50, 5.35, 5.55, 5.25, 5.40, 5.55, 6.25)
7 > angles <- matrix(
8 + cbind(C1, C2),
9 + nrow = 22,
10 + byrow = FALSE,
11 + dimnames = list(1:22, c("Narrow", "Wide")))
12 + )
13 > friedman.test(angles)
14
15      Friedman rank sum test
16
17 data: angles
18 Friedman chi-squared = 6.5455, df = 1, p-value = 0.01052
```

Chapter 14. Nonparametric Statistics

§ 14.1 Introduction

§ 14.2 The Sign Test

§ 14.3 Wilcoxon Tests

§ 14.4 The Kruskal-Wallis Test

§ 14.5 The Friedman Test

§ 14.6 Testing for Randomness

Whether the sample are random at all?

E.g. Whether the number of successful strikes are random? $\alpha = 0.05$.

Year	Number of Strikes	% Successful, y_i
1881	451	61
1882	454	53
1883	478	58
1884	443	51
1885	645	52
1886	1432	34
1887	1436	45
1888	906	52
1889	1075	46
1890	1833	52
1891	1717	37
1892	1298	39
1893	1305	50
1894	1349	38
1895	1215	55
1896	1026	59
1897	1078	57
1898	1056	64
1899	1797	73
1900	1779	46
1901	2924	48
1902	3161	47
1903	3494	40
1904	2307	35
1905	2077	40

Sol. Compute the run-up and run-down:

Year	Number of Strikes	% Successful, y_i	$\text{sgn}(y_i - y_{i-1})$	
1881	451	61	1 \rightarrow	-
1882	454	53	2 \rightarrow	+
1883	478	58	3 \rightarrow	-
1884	443	51	4 \rightarrow	+
1885	645	52	5 \rightarrow	-
1886	1432	34	6 \rightarrow	+
1887	1436	45		+
1888	906	52	7 \rightarrow	-
1889	1075	46	8 \rightarrow	+
1890	1833	52	9 \rightarrow	-
1891	1717	37	10 \rightarrow	+
1892	1298	39		+
1893	1305	50	11 \rightarrow	-
1894	1349	38	12 \rightarrow	+
1895	1215	55		+
1896	1026	59	13 \rightarrow	-
1897	1078	57	14 \rightarrow	+
1898	1056	64		+
1899	1797	73	15 \rightarrow	-
1900	1779	46	16 \rightarrow	+
1901	2924	48	17 \rightarrow	-
1902	3161	47		-
1903	3494	40		-
1904	2307	35	18 \rightarrow	+
1905	2077	40		

$w = 18$

Theorem Let W be the number of runs up and down in a sequence of $n \geq 2$ observations.

If the sequence is random, then

$$\mathbb{E}(W) = \frac{2n-1}{3} \quad \text{and} \quad \text{Var}(W) = \frac{16n-29}{90}.$$

Moreover, when n is large, namely, $n \geq 20$, then

$$\frac{W - \mathbb{E}(W)}{\sqrt{\text{Var}(W)}} = \frac{W - [2n-1]/3}{\sqrt{[16n-29]/90}} \underset{\sim}{\text{approx}} N(0, 1).$$

Sol. (Continued) $n = 25$, $w = 18$

$$\mathbb{E}(W) = \frac{2 \times 25 - 1}{3} = 16.3$$

and

$$\text{Var}(W) = \frac{16 \times 25 - 29}{90} = 4.12.$$

Hence, the z-score is

$$z = \frac{18 - 16.3}{\sqrt{4.12}} = 0.84.$$

The critical region is

$$C = \{z : |z| \geq z_{\alpha/2} = z_{0.025} = 1.96\}$$

The p -value is

$$2 \times \mathbb{P}(Z > 0.84) = 0.4009084$$

Conclusion: Fail to reject.



R code:

```
1  
2 library("snpar")  
3 y <- c(0,1,0,1,0,1,1,0,1,0,1,1,  
4        0,1,1,0,1,1,0,1,0,0,0,1)  
5 runs.test(y, exact = FALSE)  
6 runs.test(y, exact = TRUE)
```

Output:

```
1 > runs.test(y, exact = FALSE)  
2  
3      Approximate runs test  
4  
5 data: y  
6 Runs = 18, p-value = 0.03256  
7 alternative hypothesis: two.sided  
8  
9 > runs.test(y, exact = TRUE)  
10  
11      Exact runs test  
12  
13 data: y  
14 Runs = 18, p-value = 0.01624  
15 alternative hypothesis: two.sided
```

Remark The procedure that we learnt is an approximation. There is a big discrepancy for the above two p -values: one that we obtained through formula and one that is obtained by the `r` function.

Thanks for learning statistics
with me through the
semester !