Math 362: Mathematical Statistics II

Le Chen le.chen@emory.edu

Emory University Atlanta, GA

Last updated on April 24, 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

Plan

§ 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

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

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

 $1 : X \sim Binomial(n, 1/2).$ Moreover, if *m* is large, by CEE.

 $\frac{X - 2[X]}{\sqrt{Var(X)}} = \frac{X - \frac{3}{2}}{\sqrt{Nar(X)}} = \frac{N(0, 1)}{\sqrt{N(1)}}$

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

► For a random sample of size n is taken from $f_Y(y)$, in order to test $H_0: \tilde{\mu} = \tilde{\mu}_0$ vs $H_0: \tilde{\mu} \neq \tilde{\mu}_0$,

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

 \downarrow

 $(1, X \sim Binomial(n, 1/2))$ 2. Moreover, if *this* large, by CET,

 $\frac{X - 2[X]}{\sqrt{2\pi(X)}} = \frac{X - \frac{1}{2}}{\sqrt{2\pi(X)}} = N(0, 1)$

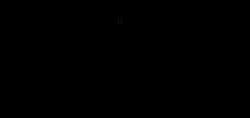
$$\mathbb{P}(\mathbf{Y} \leq \widetilde{\mu}) = \mathbb{P}(\mathbf{Y} \geq \widetilde{\mu}) = \frac{1}{2}.$$

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

$$H_0: \widetilde{\mu} = \widetilde{\mu}_0 \quad \mathrm{vs} \quad H_0: \widetilde{\mu} \neq \widetilde{\mu}_0,$$

let

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



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

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

$$H_0: \widetilde{\mu} = \widetilde{\mu}_0 \quad \mathrm{vs} \quad H_0: \widetilde{\mu}
eq \widetilde{\mu}_0,$$

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

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

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

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

$$H_0: \widetilde{\mu} = \widetilde{\mu}_0 \quad \mathrm{vs} \quad H_0: \widetilde{\mu} \neq \widetilde{\mu}_0,$$

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

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

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

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

$$H_0: \widetilde{\mu} = \widetilde{\mu}_0 \quad \mathrm{vs} \quad H_0: \widetilde{\mu}
eq \widetilde{\mu}_0,$$

let

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

\downarrow

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

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

Sign test for median of a single sample

▶ When sample size n is large:

 \blacktriangleright When sample size n is small: use the exact distribution of binomial distribution.

Sign test for median of a single sample

• When sample size n is large:

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

- **a.** 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 $z \ge z_{\alpha}$.
- **b.** 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 $z \leq -z_{\alpha}$.
- **c**. 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

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

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

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

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

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		

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

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





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

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	$_{\rm BV}$	6.95
PM	7.20 1	WR	7.35
JG	7.41	HT	7.36
DH	7.77 \angle	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 $\boldsymbol{X})$

Subject	Synovial Fluid pH	Subject	Synovial Fluid pH
HW AD	7.02 7.35	BG GL	7.34 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	$_{\rm BV}$	6.95
PM	7.20	WR	7.35
JG	7.20 1 7.41 2	HT	7.36
DH	7.77 2	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 $\boldsymbol{X})$

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

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

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 4
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	$\frac{7.41}{7.77}$ $\frac{1}{2}$	HT	7.36
DH	7.77 2	ND	6.60
ER	7.12 3	SJ	7.29
DP	7.45	BA	7.31
FF	7.28		

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

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

$$(-\infty, -z_{\alpha/2}) \cup (z_{\alpha/2}, \infty)$$
$$\parallel$$
$$(-\infty, -2.58) \cup (2.58, \infty),$$

we reject the hypothesis.

Or equivalently, the p-value is

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

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

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

$$(-\infty, -z_{\alpha/2}) \cup (z_{\alpha/2}, \infty)$$
$$||$$
$$(-\infty, -2.58) \cup (2.58, \infty),$$

we reject the hypothesis.

Or equivalently, the p-value is

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

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

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

$$(-\infty, -z_{\alpha/2}) \cup (z_{\alpha/2}, \infty)$$
$$\parallel$$
$$(-\infty, -2.58) \cup (2.58, \infty),$$

we reject the hypothesis.

Or equivalently, the p-value is

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

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

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

$$(-\infty, -z_{\alpha/2}) \cup (z_{\alpha/2}, \infty)$$
$$||$$
$$(-\infty, -2.58) \cup (2.58, \infty),$$

we reject the hypothesis.

Or equivalently, the p-value is

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

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

The exact p-value should be

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

which is smaller than $\alpha = 0.10$.

Hence, rejection!

| > pbinom(5,43,0.5) * 22 [1] 2.49951e-07 Sol 2. We can also carry out the exact computation thanks to computer: The exact p-value should be

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

which is smaller than $\alpha = 0.10$.

Hence, rejection!

| > pbinom(5,43,0.5) * 22 [1] 2.49951e-07 Sol 2. We can also carry out the exact computation thanks to computer: The exact p-value should be

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

which is smaller than $\alpha = 0.10$.

Hence, rejection!

| > pbinom(5,43,0.5) * 22 [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:

Test at $\alpha = 0.10$ that

 H_0 : consumers do not prefer B over A

vs.

 H_1 : consumers do prefer B over A.

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
В	8
А	1
No preference	1

Test at $\alpha = 0.10$ that

 H_0 : consumers do not prefer B over A

 H_1 : consumers do prefer B over A.

VS.

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
В	8
А	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.

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:

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

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:

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

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:

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

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:

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