

①

Distribution-Free Tests (in Chapter 14)

All the methods we have seen so far rely on assumption on the true distribution of data. Namely, we always assumed that $X_1, \dots, X_m \sim f_\theta$, for some $\theta \in \mathbb{R}^d$ and a known form of f_θ unknown

Ex: T-tests for means of continuous data assume $X_i \sim N(\mu, \sigma^2)$
(here $\theta = (\mu, \sigma^2)$)

Such an assumption is especially important to derive the distribution of the test statistic under the null, and hence design the testing procedure.

Parametric methods: Assume the data arise from a distribution described by a few parameters (think of $\theta = (\mu, \sigma^2)$ for Gaussians)

Nonparametric methods: Do not make parametric assumptions (most often based on ranks as opposed to raw values)

→ In this chapter, we discuss non-parametric alternatives to the one- and two-sample T-tests.

②

Examples (of when the parametric T-test goes wrong)

Extreme Outliers

Test 1: T-test comparing the two datasets

$$X = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$$

$$Y = \{7, 8, \dots, 20\}$$

$$\bar{X} = 5.5, \bar{Y} = 13.5, p\text{-value} = 1.9 \cdot 10^{-5}$$

Test 2

$$X = \{1, 2, \dots, 10\}$$

Outlier

$$Y = \{7, 8, \dots, 20, \underline{\underline{200}}\}$$

$$\bar{X} = 5.5, \bar{Y} = 25.9, p\text{-value} = 0.12$$

Why is this happening? Technically speaking, because $S_y^2 = 2335$
(while $S_x^2 = 9.2, S_y^2 = 175$)

Skew: Similarly, one can make up examples showing the distributional skew has T-tests fail.

Can you make your own?

(3)

Rk: (When to use nonparametric methods)

- With correct assumptions (e.g. normal distribution), parametric methods will be more efficient (= powerful) than nonparametric ones, but often not as much as you might think.
- If the normality assumption is grossly violated, nonparametric tests can be much more efficient and powerful than the parametric one.
- Nonparametric methods provide a well-founded way to deal with circumstances in which parametric methods perform poorly.

14.2 The sign test

Idea:

- ① T-tests are about means, and fail with outliers or skewed data
- ② The median is a much better measure of centrality, when describing data which is skewed, or with outliers
- ③ let's focus on tests about medians in such cases!

Recall: The median of a real random variable is the value $\bar{\mu}_x$ such that $P(X \leq \bar{\mu}_x) = P(X > \bar{\mu}_x) = \frac{1}{2}$

(4)

Thm: (Sign Test)

Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f_x$ be a sample from a continuous distribution with density f_x . If $n \geq 10$ and $\bar{\mu}_x = \bar{\mu}_0$, then

$$T = \frac{K - \frac{n}{2}}{\sqrt{\frac{n}{4}}} \stackrel{\text{Approx.}}{\sim} N(0, 1),$$

where $K = \sum_{i=1}^n \mathbb{1}_{\{X_i \geq \bar{\mu}_0\}}$ is the number of sample points that are bigger or equal to $\bar{\mu}_0$.

Proof: Let $Y_i = \mathbb{1}_{\{X_i \geq \bar{\mu}_0\}}$. If $\bar{\mu}_x = \bar{\mu}_0$, then $\sim \text{Binomial}(n, p)$,

where $p = P(X_i \geq \bar{\mu}_0) = \frac{1}{2}$. Conclude using the CLT.

Rk: As a corollary, we can test $H_0: \bar{\mu}_x = \bar{\mu}_0$ against

$H_1: \bar{\mu}_x > \bar{\mu}_0$, with rejection region $\{T \geq z_{\alpha}\}$

$H_1: \bar{\mu}_x < \bar{\mu}_0 \quad \Leftrightarrow \quad \{T \leq -z_{\alpha}\}$

$H_1: \bar{\mu}_x \neq \bar{\mu}_0 \quad \Leftrightarrow \quad \{|T| \geq z_{\alpha/2}\}$

(5)

Question: when are we able to use this test for $H_0: \mu_x = \mu_0$?

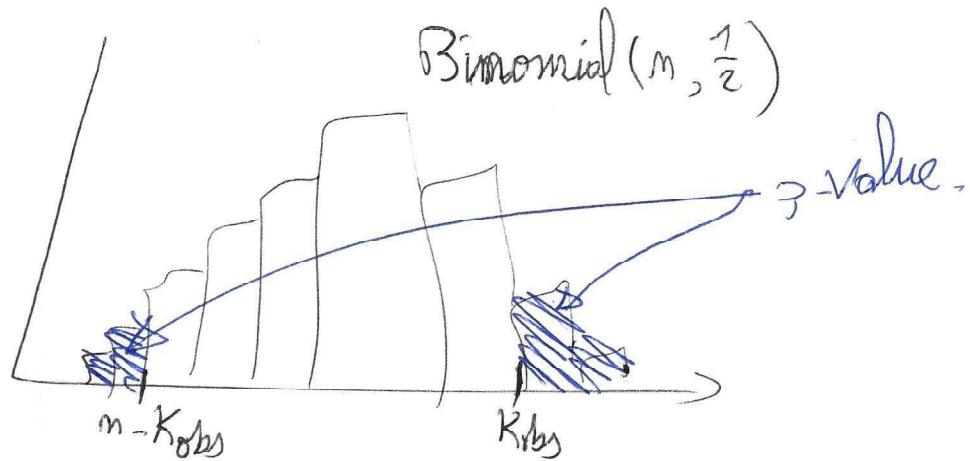
→ When f_x is symmetric about μ_x . Indeed, in such a situation,
 $\mu_x = \bar{\mu}_x$. (Prove this!)

Rk: what to do if $n < 10$?

The best way to proceed is to derive p-values directly on K with an exact test. Indeed, under H_0 , $K \sim \text{Binomial}(n, \frac{1}{2})$, so we can compute the exact p-value.

For instance, for $H_1: \bar{\mu}_x \neq \mu_0$, we get

$$\text{p-value} = P_{H_1} (K \geq K_{\text{observed}} \text{ or } K \leq n - K_{\text{observed}})$$



⑥ The sign test can also be used for paired data, by using the usual "difference" trick.

Then: let (X_1, \dots, X_n) and (Y_1, \dots, Y_n) be two paired samples.

If $\bar{\mu}_{X-Y} = 0$, then

$$K \begin{cases} \approx N\left(\frac{n}{2}, \frac{n}{4}\right) & \text{if } n \rightarrow \infty \\ \sim \text{Binomial}(n, \frac{1}{2}) & \text{if } n \leq 10 \end{cases}$$

where $K = \sum \mathbb{I}_{X_i - Y_i \geq 0}$ is the number of sample differences that are bigger or equal to 0.

Proof: Apply the previous theorem to $D_i = X_i - Y_i$.

7 A general remark about testing quantiles

Going from parametrics to nonparametrics, we basically have replaced the parameter of interest that describes a notion of centrality of the unknown distribution.

Namely, we have replaced $\mu_x = E(X)$ by $\bar{\mu}_x = \text{Median}(X)$, so that we test $H_0: \text{Median}(X) = \bar{\mu}_0$, which translates to

$$H_0: P(X \leq \bar{\mu}_0) = \frac{1}{2} .$$

This suggests to generalise this idea to any other quantile! Say you want to test

$$H_0: P(X \leq \bar{\mu}_0) = p_0$$

\uparrow

= "The p th quantile of X is equal to p ".

What would the test statistic be?

→ Still taking $K = \sum_{i=1}^n \mathbb{1}_{\{X_i > \bar{\mu}_0\}}$, we know that under the null,

$$K \sim \text{Binomial}(n, p_0) \underset{\uparrow}{\approx} N(m p_0, m p_0(1-p_0))$$

when $m p_0, m(1-p_0) \geq 5$

(8)

As a consequence, one should consider

$$T = \frac{K - mp_0}{\sqrt{mp_0(1-p_0)}},$$

and proceed as for the median (case $p = \frac{1}{2}$)

Supplementary Question: How do we do if $mp_0 < 5$?

Two possibilities: 1) If m is small, use the exact binomial model (and p -value)

2) If m is large, use the Poisson approximation

$$\text{Binomial}(m, p_0) \approx \text{Poisson}(mp_0)$$

9

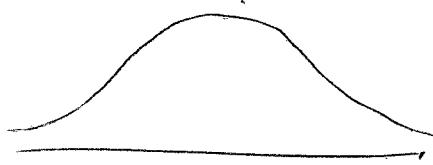
14.3 Wilcoxon Tests

We now move to a more elaborate family of methods based on ranks, but that require a bit more of hypotheses.

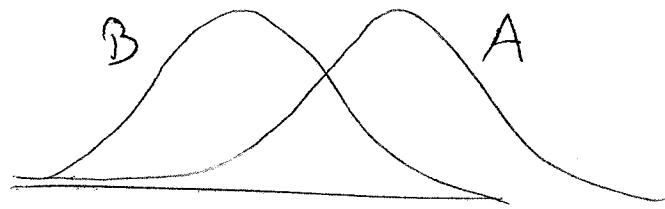
- Ex:
- Suppose we have m and n experimental units
 - $\{ \cdot \}_{m \text{ control}}$
 - $\{ \cdot \}_{n \text{ treatment}}$
 - The assignments is made at random
 - We're interested in testing if the treatment does has positive effect or not.

Very often, a positive effect would translate in a shift on the distribution,

$$H_0: A = B$$



$$H_1: A > B$$



without a change of shape, whatever the shape

→ let's use the shape shift phenomenon in a test.

- Idea:
- Group all m and n observations together and rank them in order of increasing size
 - Calculate the sum of ranks of the control group
 - If this sum is too large, reject the null,

(10)

Thm (Wilcoxon Signed Ranks Test)

let (X_1, \dots, X_n) be independent observations drawn from p.d.f's f_{X_1}, \dots, f_{X_n} all of which are continuous and symmetric with equal mean μ .

Then for testing

$$H_0: \mu = \mu_0 \quad (\text{vs}) \quad H_1: \mu \neq \mu_0$$

We use

$$T = \sum_{i=1}^n R_i Z_i \quad \begin{cases} \sim N\left(\frac{n(n+1)}{4}, \frac{n(n+1)(2n+1)}{24}\right) & \text{if } n \rightarrow \infty \\ \sim P(T=t) = \left(\frac{1}{2}\right)^n C(t) & \text{if } n \leq 12 \end{cases}$$

where the rank of $|X - \mu_0|$, and

$$Z_i = \begin{cases} 1 & \text{if } X_i - \mu_0 < 0 \\ 0 & \text{if } X_i - \mu_0 \geq 0 \end{cases}$$

and

$C(t)$ is the coefficient of e^{tx} in the expansion of $\prod_{i=1}^n (1 + e^{ix})$

(11)

- Rk:
- Again, we did not need to assume a specific family of distributions
 - The X_i 's don't need to be identically distributed
 - . We have both asymptotic and nonasymptotic

Proof: If H_0 is true, since $\bar{\mu}_{X_i} = \mu_{X_i}$ because of symmetry, T has the same distribution as

$$U = \sum_{i=1}^n U_i, \text{ where } \begin{cases} P(U_i=0) = \frac{1}{2} \\ P(U_i=1) = \frac{1}{2} \end{cases}, U_i \text{'s } \perp \!\! \perp$$

Hence, we can do the reasoning on U . To describe its distribution, we compute its moment generating function

$$\begin{aligned} M_U(t) &= E\left(e^{\sum_{i=1}^n U_i t}\right) \\ &= \prod_{i=1}^n E(e^{U_i t}) \\ &= \prod_{i=1}^n \left(\frac{1}{2} e^{0t} + \frac{1}{2} e^{1 \cdot t}\right) \\ &= \left(\frac{1}{2}\right)^n \prod_{i=1}^n (1 + e^{it}) \end{aligned}$$

(12)

We conclude by noticing that since U is discrete,

$$\begin{aligned} M_U(t) &= \mathbb{E}(e^{tU}) \\ &= \sum_{x \in \mathbb{N}} e^{tx} P(U=x), \end{aligned}$$

and by developing the previous product]

Rk: - The asymptotic normality is obtained using a version of the CLT (on the U_i 's) that is out of the scope of this course. If interested, see "CLT for triangular arrays".

The following result is especially designed for two-sample problems, as opposed to the previous one.

(13)

Thm (Wilcoxon Ranks Sum Test / Mann-Whitney)

Let (X_1, \dots, X_m) and (Y_1, \dots, Y_n) be two independent samples from p.d.f.'s f_X and f_Y which are continuous, symmetric, of the same shape and the same standard deviation. (\equiv one might only be the shifted version of the other)

If $\mu_X = \mu_Y$,

$$T = \sum_{i=1}^{m+n} R_i Z_i \sim N\left(\frac{m m}{2}, \frac{m m (m+n+1)}{12}\right)$$

if $m, n \geq 10$,

where R_i is the rank of the i^{th} observation in the joint sample

$(X_1, \dots, X_m, Y_1, \dots, Y_n)$

$$Z_i = \begin{cases} 1 & \text{if the } i^{\text{th}} \text{ observation comes from } f_X \\ 0 & \text{if the } i^{\text{th}} \text{ observation comes from } f_Y \end{cases}$$

Proof: Similar