

Lecture 10: Hypothesis Testing

Ma Xuejun

School of Mathematical Sciences

Soochow University

<https://xuejunma.github.io>



Hypothesis Testing

For a hypothesis testing problem:

- Underlying truth: H_0 is true or H_1 is true
- Goal: sufficient evidence to reject H_0 ?
- Action: reject H_0 or not reject H_0

	Decision	
	Retain H_0	Reject H_0
H_0 is true	✓	Type I error(false positive)
H_1 is true	Type II error(false negative)	✓

- Without sufficient evidence, we do not reject H_0 . It does not mean we believe it is correct
- Obviously, the setting prefers H_0

Evaluation of a test

- With a test, we hope we can do correct justifications.
- It means minimizing the Type I error and Type II error.
- For the type II error, we define the power function.

Definition: Power function

Suppose we reject H_0 when $T(X_1, \dots, X_n) \in R$. The power function is defined as

$$\beta(\theta) = P(T(X_1, \dots, X_n) \in R | \theta).$$

Remark.

- The power function is a function about θ
- When $\theta \in \Theta_1$, it measures the probability that the test correctly rejects H_0
- When $\theta \in \Theta_0$, it measures the type I error.

Evaluation

The general strategy to construct a test is

- (1) Fixe $\alpha \in [0, 1]$
- (2) Try to maximize $\beta(\theta)$ for $\theta \in \Theta_0$, subject to $\beta(\theta) \leq \alpha$ for $\theta \in \Theta_0$

Example. $X_1, \dots, X_n \sim N(\theta, \sigma^2)$ with σ^2 known. Suppose we test

$$H_0 : \theta = \theta_0 \quad \text{versus} \quad H_1 : \theta > \theta_0.$$

This is called a one – sided alternative. Suppose we reject H_0 if $T_n > c$ where

$$T_n = \frac{\bar{X}_n - \theta_0}{\sigma/\sqrt{n}}.$$

Then, the power function is

$$\begin{aligned} \beta(\theta) &= P\left(\frac{\bar{X}_n - \theta_0}{\sigma/\sqrt{n}} > c; \theta\right) = P\left(\frac{\bar{X}_n - \theta}{\sigma/\sqrt{n}} > c + \frac{\theta_0 - \theta}{\sigma/\sqrt{n}}; \theta\right) \\ &= P\left(Z > c + \frac{\theta_0 - \theta}{\sigma/\sqrt{n}}\right) = 1 - \Phi\left(c + \frac{\theta_0 - \theta}{\sigma/\sqrt{n}}\right), \end{aligned}$$

where $Z \sim N(0, 1)$ and Φ is the CDF for Z . Now,

$$\sup_{\theta \in \Theta_0} \beta(\theta) = \beta(\theta_0) = 1 - \Phi(c).$$

To get a *size* α test, set $1 - \Phi(c) = \alpha$ so that

$$c = z_\alpha = \Phi^{-1}(1 - \alpha).$$

Our test is to reject H_0 when $T_n = \frac{\bar{X}_n - \theta_0}{\sigma/\sqrt{n}} > z_\alpha$.

Now, let's consider the **two-sided alternative**, that

$$H_0 : \theta = \theta_0 \quad \text{versus} \quad H_1 : \theta \neq \theta_0.$$

We will reject H_0 if $|T_n| > c$. The power function is

$$\begin{aligned}\beta(\theta) &= P(T_n < -c; \theta) + P(T_n > c; \theta) \\ &= P\left(\frac{\bar{X}_n - \theta_0}{\sigma/\sqrt{n}} < -c; \theta\right) + P\left(\frac{\bar{X}_n - \theta_0}{\sigma/\sqrt{n}} > c; \theta\right) \\ &= \Phi\left(-c + \frac{\theta_0 - \theta}{\sigma/\sqrt{n}}\right) + 1 - \Phi\left(c + \frac{\theta_0 - \theta}{\sigma/\sqrt{n}}\right) \\ &= \Phi\left(-c + \frac{\theta_0 - \theta}{\sigma/\sqrt{n}}\right) + \Phi\left(-c - \frac{\theta_0 - \theta}{\sigma/\sqrt{n}}\right)\end{aligned}$$

The size is $\beta(\theta_0) = 2\Phi(-c)$. Let it equal to α , then

$c = -\Phi^{-1}(1 - \alpha/2) = z_{\alpha/2}$. The test is to reject H_0 when $|\frac{\bar{X}_n - \theta_0}{\sigma/\sqrt{n}}| > z_{\alpha/2}$

- Neyman-Pearson Test
- Wald Test
- Likelihood Ratio Test (LRT)
- Score Test
- Permutation Test
- Bootstrap Test

Now we discuss them one by one.

The Neyman-Pearson Test

- The Neyman-Pearson test considers only simple null and simple alternative, which means the test

$$H_0 : \theta = \theta_0 \quad \text{versus} \quad H_1 : \theta = \theta_1.$$

Definition: Neyman-Pearson Test

Let $L(\theta) = f(X_1, \dots, X_n; \theta)$ and

$$T_n = L(\theta_1)/L(\theta_0).$$

Suppose we reject H_0 if $T_n > k$ where k is chosen so that

$$P\left(T(X_1, \dots, X_n) > k; \theta = \theta_0\right) = \alpha,$$

then it is called a Neyman-Pearson Test.

- The test statistic is the ratio of two joint densities. It is to check with which likelihood, the data is more possible.
- It is quite limited, since it requires both the null and the alternative are simple.

The Neyman-Pearson Test

Definition 8.3.11 : Uniformly Powerful Tests

Let C_α be a collection of level α for $H_0 : \theta \in \Theta_0$ and $H_1 : \theta \in \Theta_1$. A test in C_α with power function $\beta(\theta)$ is uniformly most powerful (UMP) if for every $\beta'(\theta)$ which is the power function of any other test in C_α , then

$$\beta(\theta) \geq \beta'(\theta), \quad \theta \in \Theta_1.$$

- The requirements in Definition 8.3.11 are so strong that UMP tests may not exist in realistic problem.
- In the simple null and simple alternative case, it exists, which is the Neyman-Pearson test.

- **Theorem 8.3.12 (Neyman-Pearson Lemma)** Consider testing $H_0 : \theta = \theta_0$ v.s. $H_1 : \theta = \theta_1$, where the pdf or pmf corresponding to θ_i is $f(\mathbf{x}|\theta_i)$, $i = 0, 1$, using a test with rejection region R that satisfies

$$\mathbf{x} \in R \text{ if } f(\mathbf{x}|\theta_1) > k f(\mathbf{x}|\theta_0) \quad (2.1)$$

and

$$\mathbf{x} \in R^c \text{ if } f(\mathbf{x}|\theta_1) < k f(\mathbf{x}|\theta_0) \quad (2.2)$$

for some $k \geq 0$ and

$$\alpha = P_{\theta_0}(\mathbf{X} \in \mathbf{R})$$

Then

- (*Sufficiency*) Any test that satisfies (2.1) and (2.2) is a UMP level α test,
- (*Necessity*) If there exists a test that satisfying (2.1) and (2.2) with $k > 0$, then every UMP level α test is a size α test (satisfies (2.2)) and every UMP level α test satisfies (2.2) except perhaps on a set \mathbf{A} satisfying $P_{\theta_0}(\mathbf{X} \in \mathbf{A}) = P_{\theta_1}(\mathbf{X} \in \mathbf{A}) = 0$.

- **Corollary 8.3.13** Suppose that $T(\mathbf{x})$ is a sufficient statistic for θ and $g(t|\theta_i)$ is the pdf or pmf of T corresponding to θ_i , $i = 0, 1$. Then any test based on T with rejection region S (a subset of the sample space of T) is a UMP level α test if it satisfies

$$t \in S \text{ if } g(t|\theta_1) > kg(t|\theta_0)$$

and

$$t \in S^c \text{ if } g(t|\theta_1) < kg(t|\theta_0)$$

for some k , where

$$\alpha = P_{\theta_0}(T \in S) \quad (2.3)$$

Example 8.3.14 (UMP binomial test) Let $X \sim \text{Binomial}(2, \theta)$. We want to test $H_0 : \theta = \frac{1}{2}$, versus $H_1 : \theta = \frac{3}{4}$. Calculating the ratios of the pmfs gives

$$\frac{f(0|\theta = \frac{1}{2})}{f(0|\theta = \frac{3}{4})} = \frac{1}{4}, \quad \frac{f(1|\theta = \frac{1}{2})}{f(1|\theta = \frac{3}{4})} = \frac{3}{4}, \quad \text{and} \quad \frac{f(2|\theta = \frac{1}{2})}{f(2|\theta = \frac{3}{4})} = \frac{9}{4}.$$

- $\frac{3}{4} < k < \frac{9}{4}$, the Neyman-Pearson Lemma says that the test that rejects H_0 if $X = 2$ is the UMP level $\alpha = P(X = 1|\theta = \frac{1}{2}) = \frac{1}{4}$ test.
- $\frac{1}{4} < k < \frac{3}{4}$, the Neyman-Pearson Lemma says that the test that rejects H_0 if $X=1$ or 2 is the UMP level $\alpha = P(X = 1 \text{ or } 2|\theta = \frac{1}{2}) = \frac{3}{4}$ test.

...

- **Example 8.3.15 (UMP normal test)** Let X_1, \dots, X_n be a random sample from a $N(\theta, \sigma^2)$ population, σ^2 known.
- The sample mean \bar{X} is a sufficient statistic for θ .
- Consider testing $H_0 : \theta = \theta_0$, versus $H_1 : \theta = \theta_1$, where $\theta_0 > \theta_1$.
- The inequality $g(\bar{x}|\theta_1) > kg(\bar{x}|\theta_0)$, is equivalent to

$$\bar{x} < \frac{(2\sigma^2 \log k)/n - \theta_0^2 + \theta_1^2}{2(\theta_0 - \theta_1)}.$$

The fact that $\theta_1 - \theta_0 < 0$ was used to obtain the inequality.

- The right-hand side increases from $-\infty$ to ∞ as k increases from 0 to ∞ .
- Thus, by Corollary 8.3.13, the test with rejection region $\bar{x} < c$ is the UMP level α test, where $\alpha = P_{\theta_0}(\bar{X} < c)$.
- If a particular α is specified, then the UMP test rejects H_0 if $\bar{X} < c = -\sigma z_\alpha / \sqrt{n} + \theta_0$. This choice of c ensures that (2.3) is true.
- In the case of composite hypothesis, the UMP test can be derived with the Neyman-Pearson Lemma.

The Wald Test

- Assume there is an asymptotic normal estimator $\hat{\theta}_n$, where $\hat{\theta}_n - \theta \xrightarrow{d} N(0, \sigma_n^2)$
- If $H_0 : \theta = \theta_0$ is true, then there is $\hat{\theta}_n - \theta \xrightarrow{d} N(0, \sigma_n^2)$
- we can construct a test statistic

$$T_n = \frac{\hat{\theta}_n - \theta_0}{\hat{\sigma}_n}$$

- If H_0 is true, $T_n \xrightarrow{d} N(0, 1)$, which concentrates at 0. So, if T_n is too large/small, we reject H_0 .
- This kind of test is called the Wald Test.

Example.

- With Bernoulli data, to test $H_0 : p = p_0$ and $H_1 : p \neq p_0$, recall that $\sqrt{n}(\bar{X} - p) \xrightarrow{d} N(0, p(1-p))$, we can construct a Wald test

$$T_n = \left| \frac{\bar{X} - p_0}{\sqrt{\hat{p}(1-\hat{p})/n}} \right| > c,$$

where $c = \Phi^{-1}(1 - \alpha/2) = z_{\alpha/2}$.

The Wald Test

- Consider MLE $\hat{\theta}_n$. According to the asymptotic normality of MLE, there is

$$\sqrt{n} \frac{\hat{\theta}_n - \theta}{\sqrt{1/I(\theta)}} \xrightarrow{d} N(0, 1).$$

So we can construct a test w.r.t. MLE, which is to reject null hypothesis when

$$T_n = \frac{\hat{\theta}_n - \theta_0}{\sqrt{1/nI(\hat{\theta})}} > c.$$

- If it happens that \bar{X} is an estimator for θ . According to CLT, $\sqrt{n}(\bar{X} - \theta)/\sigma \xrightarrow{d} N(0, 1)$. So we can also build a Wald test based on the average
- Usually, σ_n is a function of θ . Since the truth is unknown, we can either apply θ_0 or $\hat{\theta}$ in practice.
- The Wald test requires asymptotic normality, so it works for *large* sample size only.

The Likelihood Ratio Test

- Neymann-Pearson test is the ratio of likelihoods w.r.t. two values
- For composite null and alternative, we can generalize the idea

Definition: Likelihood Ratio Test (LRT)

The LRT statistic for testing $H_0 : \theta \in \Theta_0$ versus $H_1 : \theta \in \Theta_1$ is

$$\lambda(x_1, \dots, x_n) = \frac{\sup_{\theta \in \Theta_0} f_{X_{1:n}}(x_1, \dots, x_n; \theta)}{\sup_{\theta \in \{\Theta_0 \cup \Theta_1\}} f_{X_{1:n}}(x_1, \dots, x_n; \theta)}$$

A LRT is any test that has a rejection region of the form $\{(x_1, \dots, x_n); \lambda(x_1, \dots, x_n) \leq c\}$ for any constant $c \in [0, 1]$.

- Θ_0 : null parameter space; Θ : the whole parameter space
- According to the definition of MLE, the LRT statistic can be written as

$$\lambda(x_1, \dots, x_n) = \frac{L(\hat{\theta}_0)}{L(\hat{\theta})}$$

- If it is small, then we reject H_0 .

LRT: Example

Example. Suppose that $X_i \stackrel{i.i.d}{\sim} N(\theta, 1)$ and suppose we want to test $H_0 : \theta = \theta_0$ versus $H_1 : \theta \neq \theta_0$.

Recall that the MLE is $\hat{\theta} = \bar{X}_n$. So the LRT statistic is

$$\lambda(x_1, \dots, x_n) = \frac{\frac{1}{(\sqrt{2\pi})^n} e^{-\frac{\sum (x_i - \theta_0)^2}{2}}}{\frac{1}{(\sqrt{2\pi})^n} e^{-\frac{\sum (x_i - \bar{x})^2}{2}}} = \frac{\exp\left\{-\frac{\sum (x_i - \theta_0)^2}{2}\right\}}{\exp\left\{-\frac{\sum (x_i - \bar{x})^2}{2}\right\}}.$$

Since $\sum (x_i - \theta_0)^2 = \sum (x_i - \bar{x})^2 + n \sum (\bar{x} - \theta_0)^2$, we have

$\lambda(x_1, \dots, x_n) = \exp\left\{-\frac{n}{2}(\bar{x} - \theta_0)^2\right\}$. Since it is monotone with $|\bar{x} - \theta_0|$, so the rejection region is equivalent with

$$\{x \in \mathbb{R}^n : |\bar{x} - \theta_0| \geq c\}.$$

Since $\bar{x} - \theta_0 \sim N(0, 1/n)$ under null hypothesis, the level of the test is $2\Phi(-\sqrt{nc})$. For a level α test, we have $c = \Phi^{-1}(1 - \alpha/2)/\sqrt{n}$.

- Example 8.2.3 (Exponential LRT)** X_1, X_2, \dots, X_n be a random

sample from $f(x|\theta) = \begin{cases} e^{-(x-\theta)}, & x \geq \theta, \\ 0, & x < \theta. \end{cases}$

$$L(\theta|\mathbf{x}) = \begin{cases} e^{-\sum x_i + n\theta}, & \theta \leq x_{(1)}, \\ 0, & \theta > x_{(1)}. \end{cases}$$

$H_0 : \theta \leq \theta_0$, versus $H_1 : \theta > \theta_0$

$$\sup_{\theta \leq \theta_0} L(\theta|\mathbf{x}) = \begin{cases} e^{-\sum x_i + n\theta_0}, & \theta < x_{(1)}, \\ e^{-\sum x_i + nx_{(1)}}, & \theta \geq x_{(1)}. \end{cases}$$

$$\sup_{\theta} L(\theta|\mathbf{x}) = e^{-\sum x_i + nx_{(1)}}.$$

Therefore,

$$\lambda(\mathbf{x}) = \begin{cases} 1, & x_{(1)} \leq \theta_0, \\ e^{-n(x_{(1)} - \theta_0)}, & x_{(1)} > \theta_0. \end{cases}$$

$$\{\mathbf{x} : \lambda(\mathbf{x}) \leq c\} = \{\mathbf{x} : x_{(1)} \geq \theta_0 - \frac{\log c}{n}\}$$

Note that the rejection region depends on the sample only through $x_{(1)}$.

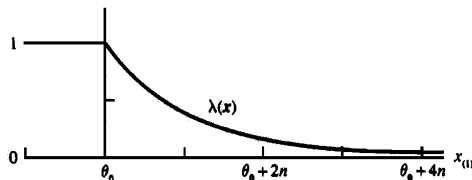


Figure 8.2.1 $\lambda(x)$: a function only of $x_{(1)}$

LRT: Theorem

- Can we always find a proper distribution for the LRT statistic?
- Not always, but asymptotically, yes.

Theorem: LRT statistics

Let $X_i \overset{i.i.d.}{\sim} F_X(\cdot; \theta^*)$ with $f_X(\cdot; \theta^*)$ as the associated PDF. Let $\hat{\theta}_n$ be the MLE. Consider the testing $H_0 : \theta = \theta_0$ versus $H_1 : \theta \neq \theta_0$ where $\theta \in \mathbb{R}$.

The under H_0 ,

$$-2 \log(\lambda(X_1, X_2, \dots, X_n)) \xrightarrow{d} \chi_1^2.$$

- Regularity conditions for MLE normality
- To construct a level α test, we can set the rejection region as $-2 \log(\lambda(X_1, X_2, \dots, X_n)) \geq \chi_{1, \alpha}^2$, where $\chi_{1, \alpha}^2$ is the $1 - \alpha$ quantile for $\chi_{1, \alpha}^2$ distribution.
- If $\theta = (\theta_1, \theta_2, \dots, \theta_k)$. Then, under the regularity conditions,

$$T_n = -2 \log(\lambda(X_1, X_2, \dots, X_n)) \xrightarrow{d} \chi_v^2, \quad v = \dim(\Theta) - \dim(\Theta_0).$$

Proof

- Under the regularity conditions, we have the Taylor expansion for the log-likelihood function $l(\theta)$ close to the point $\hat{\theta}$:

$$l(\theta) \approx l(\hat{\theta}) + l'(\hat{\theta})(\theta - \hat{\theta}) + l''(\hat{\theta}) \frac{(\theta - \hat{\theta})^2}{2} = l(\hat{\theta}) + l''(\hat{\theta}) \frac{(\theta - \hat{\theta})^2}{2}$$

- The expression for LRT statistic is

$$\begin{aligned} -2 \log(\lambda(X_1, X_2, \dots, X_n)) &= -2l(X_1, X_2, \dots, X_n; \theta_0) \\ &\quad + 2l(X_1, X_2, \dots, X_n; \hat{\theta}) \\ &\approx 2l(\hat{\theta}) - 2l(\hat{\theta}) - l''(\hat{\theta})(\theta - \hat{\theta})^2 \\ &= -l''(\hat{\theta})(\theta - \hat{\theta})^2 \\ &= \frac{-l''(\hat{\theta})}{I_n(\theta_0)} \times I_n(\theta_0)(\hat{\theta} - \theta_0)^2 = A_n \times B_n \end{aligned}$$

- Note that $A_n \xrightarrow{P} 1$ according to WLLN and $\sqrt{B_n} \xrightarrow{d} N(0, 1)$, so that $B_n \xrightarrow{d} \chi_1^2$. According to Slutsky's theorem, the result follows.

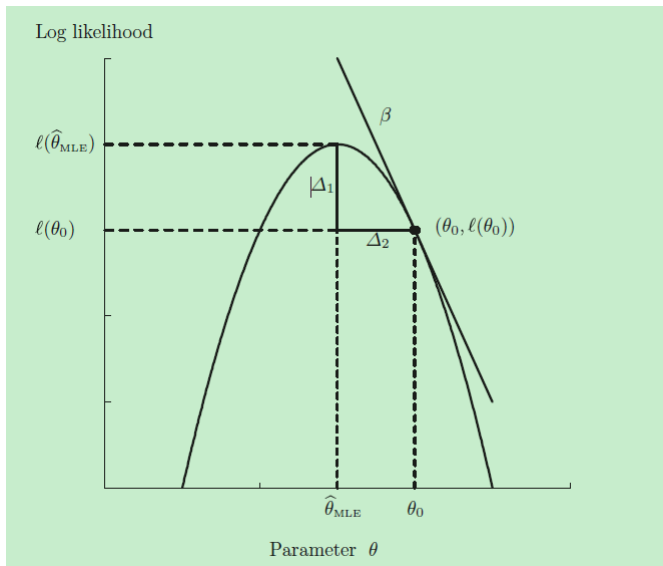
11. *Journal of the American Medical Association*, 277, 1996, 1033-1037.

- $$S^2(\theta_0) = \frac{1}{n} \sum_{i=1}^n \frac{1}{\pi(\theta_0)} \left(\frac{\partial \log \pi(\theta)}{\partial \theta} \right)^2 \bigg|_{\theta=\theta_0}$$

$$M : \quad \alpha(\rho) \quad \partial \rho(\rho) \quad \Gamma(\Gamma(\rho)) \quad \partial \alpha(\rho) \quad \Gamma(\rho) \quad \Gamma[\Gamma(\Gamma(\rho))] \quad \Gamma(\rho)$$

25 / 62

- where $\hat{\theta}^*$ lies between $\hat{\theta}_{ML}$ and θ_0 . $I_n(Y, \hat{\theta}^*) \xrightarrow{P} I(\theta_0)$



Normal model with known variance

Suppose that Y_1, \dots, Y_n iid $N(\mu, 1)$. $H_0 : \mu = \mu_0$, then

$$\ell(\mu) = \log L(\mu|Y) = -\frac{n}{2} \log 2\pi - \frac{1}{2} \sum_{i=1}^n (Y_i - \mu)^2$$

$$S(\mu) = \frac{\partial}{\partial \mu} \ell(\mu) = \sum_{i=1}^n (Y_i - \mu)$$

$$I_T(Y, \mu) = \frac{\partial}{\partial \mu} S(\mu) = n$$

So that $M\hat{L}E = \bar{Y}$, and $I_T(\mu) = E[I_T(Y, \mu)] = n$. Hence

$$T_W = \frac{(\bar{Y} - \mu_0)^2}{n^{-1}} = (\bar{Y} - \mu_0)^2$$

$$T_S = \frac{\left[\sum_{i=1}^n (Y_i - \mu) \right]^2}{n} = (\bar{Y} - \mu_0)^2$$

$$T_{LR} = -2 \left[-\frac{1}{2} \sum_{i=1}^n (Y_i - \mu_0)^2 + \frac{1}{2} \sum_{i=1}^n (Y_i - \bar{Y})^2 \right] = (\bar{Y} - \mu_0)^2$$

p-values

- Given α , we construct a level α test
- With data, we calculate the statistic and decide whether to *reject* or *retain* H_0
- If α changes, should we do all the steps again?

Definition: P-values

A p -value $p(X)$ is a test statistic with $p(X) \in [0, 1]$. Small values of p indicate that H_1 is true. A p -value is valid if for every $\theta \in \Theta_0, \alpha \in [0, 1]$,

$$P(p(X) \leq \alpha; \theta) \leq \alpha.$$

- $p(X)$ is a test statistic. With the statistic $p(X)$, the level α test is to reject H_0 when $p < \alpha$. The power function w.r.t. this test is

$$\beta(\theta) = P(p(X) \leq \alpha; \theta).$$

- Therefore, it can be **viewed as the smallest α at which we would reject H_0 .**

p-values

- Question: how to find this test statistic?

Theorem: *P*-values

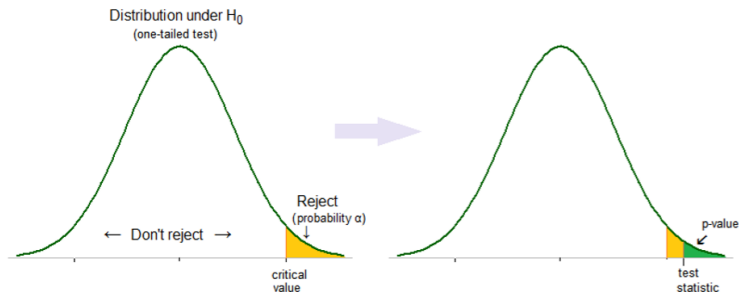
Let $W(X)$ be a test statistic such that large values of W indicate that H_1 is true. For each $x \in X$, define

$$p(x) = \sup_{\theta \in \Theta_0} P(W(X) \geq W(x); \theta),$$

then $p(X)$ is a valid *p*-value.

- Note that $W(x)$ should satisfy that reject H_0 when $T(x) > c$.
- This is the general way to find the *p*-value. We define a test statistic first, and then define *p* be the probability that the statistic is no smaller than the observation.
- p*-value may change for different test statistic, even with the same data. So when we specify *p*-value, we should specify the test statistic.

Remarks



- p -value is the probability for the test statistic under null, not the probability of H_0
- Why p -value is useful, not the test statistic?

Theorem

Under $H_0, p \sim Unif(0, 1)$

We never know what the test statistic means, but we can achieve the information in p -value quickly

Example: *p*-values

Let $X_1, \dots, X_n \stackrel{i.i.d}{\sim} N(0, 1)$. Test that

$$H_0 : \theta = \theta_0 \quad \text{versus} \quad H_1 : \theta \neq \theta_0.$$

We reject when $|T_n| = |\sqrt{n}(\bar{X}_n - \theta_0)|$ is large.

- Let t_n be the observed value of T_n . Let $Z \sim N(0, 1)$. Then,

$$p = P(|\sqrt{n}(\bar{X}_n - \theta_0)| > t_n) = P(|Z| > t_n) = 2\Phi(-|t_n|).$$

Now, we can return the *p*-value to the researcher, with which the researcher can easily tell **how strong the evidence** is to reject H_0 .

Under H_0 , the limiting distribution of T_{mn} as $m, n \rightarrow \infty$ with $m/(m+n) \rightarrow \lambda \in (0, 1)$, is normal with mean 0 and variance $\text{Kurt}(F_0) - 1$. In the introduction to Boos et al. (1989), four bootstrap resampling plans are discussed:

- I. Draw both bootstrap samples independently and with replacement from the pooled set $\{X_1, \dots, X_m, Y_1, \dots, Y_n\}$.
- II. Draw X_1^*, \dots, X_m^* with replacement from $\{X_1, \dots, X_m\}$ and independently draw Y_1^*, \dots, Y_n^* with replacement from $\{Y_1, \dots, Y_n\}$.
- III. Draw both bootstrap samples independently and with replacement from the pooled set of residuals $\{X_1 - \bar{X}, \dots, X_m - \bar{X}, Y_1 - \bar{Y}, \dots, Y_n - \bar{Y}\}$.
- IV. Draw X_1^*, \dots, X_m^* with replacement from $\{X_1/s_X, \dots, X_m/s_X\}$ and independently draw Y_1^*, \dots, Y_n^* with replacement from $\{Y_1/s_Y, \dots, Y_n/s_Y\}$.

“99 Rule”

Consider a situation where the statistic T is continuous, and a parametric bootstrap gives the exact sampling distribution as B grows large.

- T_0, T_1^*, \dots, T_B^* are iid, all $(B + 1)!$ orderings are equally likely, and p_B has a discrete uniform distribution,

$$P(p_B = 0) = P(p_B = 1/B) = \dots = P(p_B = 1) = \frac{1}{1 + B}.$$

- The test defined by the rejection region $p_B \leq \alpha$ has **exact level** α if $(B + 1)\alpha$ is an integer.
- So, for small B one should use values like $B = 19$ or 39 or 99 to get standard α levels. We call this the “99 rule” .
- The “99 rule” should be followed generally in bootstrap testing situations.

Regression Settings

We consider typical regression settings based on iid random pairs $(Y_1, \mathbf{X}_1), \dots, (Y_n, \mathbf{X}_n)$, or $(Y_1, \mathbf{x}_1), \dots, (Y_n, \mathbf{x}_n)$, where the explanatory vectors \mathbf{x}_i are viewed as fixed constants.

- In the random pairs case, it is natural to draw with replacement from the set of pairs resulting in a bootstrap resample $(Y_1^*, \mathbf{X}_1^*), \dots, (Y_n^*, \mathbf{X}_n^*)$.
- This bootstrap method is very general and applies to almost any regression method.
- The assumed model used to derive estimators does not need to be true in order for bootstrap estimates to be consistent.
- We call this method the *random pairs* bootstrap although it is really just the standard nonparametric bootstrap method.

There are a few reasons, however, to consider other bootstrap approaches in regression settings:

- Inference in regression setting is usually carried out conditional on the explanatory vectors regardless of whether they are considered fixed or random. The random pairs bootstrap, however, gives unconditional estimates.
- The random pairs bootstrap does not take advantage of any model assumptions such as an additive error structure with **homogeneous errors**. This nonparametric aspect of the random pairs bootstrap gives it strong robustness to model assumptions, but also can result in much less efficient procedures.

For these reasons, let us consider the *residual-based* bootstrap that is appropriate for additive errors models of the form $Y_i = g(\mathbf{x}_i, \beta) + e_i$, where g is a known function and e_1, \dots, e_n are iid random errors.

- Defining the residuals as $\hat{e}_i = Y_i - g(\mathbf{X}_i, \hat{\beta})$, draw bootstrap errors e_1^*, \dots, e_n^* with replacement from the set

$$\left\{ (\hat{e}_i - \bar{\hat{e}}) / \sqrt{1 - p/n}, i = 1, \dots, n \right\}.$$

- Then form the bootstrap responses $Y_i^* = g(\mathbf{X}_i, \hat{\beta}) + e_i^*, i = 1, \dots, n$.

- If the model is linear, then the least squares estimator in the bootstrap world is $\hat{\beta}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}^*$ with variance $\text{Var}^*(\hat{\beta}^*) = \hat{\sigma}^2 (\mathbf{X}^T \mathbf{X})^{-1}$, where

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n [(\hat{e}_i - \bar{\hat{e}}) / \sqrt{1 - p/n}]^2 = \frac{1}{n - p} \sum_{i=1}^n (\hat{e}_i - \bar{\hat{e}})^2.$$

- Further, if the first column of \mathbf{X} is a column of ones, then $\bar{\hat{e}} = 0$, and we recognize the bootstrap estimate of $\text{Var}(\hat{\beta})$ is the same as the usual unbiased one.

Permutation and Rank Tests

- An example of two treatment effect test based on permutation method.
- six students are divided into two groups of size 2 and 4, taught with different method, the scores of students using standard method are

$$x_1 = 6, x_2 = 8$$

the scores of new method group are

$$x_3 = 7, x_4 = 18, x_5 = 11, x_6 = 9$$

As we all know, in parametric situation, we can use t-test.

In nonparametric situation, we can also use the statistic below:

$$t(X, Y) = \frac{\bar{Y} - \bar{X}}{\sqrt{s_p^2(\frac{1}{m} + \frac{1}{n})}}$$

where

$$s_p^2 = \left\{ \sum (X_i - \bar{X})^2 + \sum (Y_i - \bar{Y})^2 \right\} / (m + n - 2)$$

If t is large, then one might be convinced that the new method is better than the standard one.

- In nonparametric situation, another common used statistic is

W = the sum of the ranks of the Y values

when both X and Y samples are thrown together and ranked from smallest to largest.

- Let Z denote the joint sample of both X and Y together:
 $Z=(X,Y)=(6,8,7,18,11,9)$,
- The ranks of these observed values are
 $(1,3,2,6,5,4)$, $W=2+6+5+4=17$.
- the new method is better, the W is expected to be large.
- t and W are statistics for our testing problem, we should find the distribution of each.
- There are $\binom{6}{2} = 15$ different ways of treating.

The possible samples and the values of t and W are listed below.

Table 12.1 All Possible Permutations for Example Data

	X Sample		Y Sample		$\sum Y_i$		t	W	
1.	6	8	7	18	11	9	45	1.17	17
2.	7	8	6	18	11	9	44	0.91	16
3.	18	8	7	6	11	9	33	-1.36	12
4.	11	8	7	18	6	9	40	0.12	13
5.	9	8	7	18	11	6	42	0.49	14
6.	6	7	8	18	11	9	46	1.47	18
7.	6	18	7	8	11	9	35	-0.84	14
8.	6	11	7	18	8	9	42	0.49	15
9.	6	9	7	18	11	9	34	-1.08	13
10.	7	18	6	8	11	9	34	-1.08	13
11.	18	11	7	6	8	9	30	-2.98	10
12.	11	9	7	18	6	8	39	-0.06	12
13.	7	11	6	18	8	9	41	0.30	14
14.	7	9	6	18	11	8	43	0.69	15
15.	18	9	7	6	11	8	32	-1.72	11

t	-2.98	-1.72	-1.36	-1.08	-0.84	-0.06	0.12
$P(t)$	$\frac{1}{15}$	$\frac{1}{15}$	$\frac{1}{15}$	$\frac{1}{15}$	$\frac{1}{15}$	$\frac{1}{15}$	$\frac{1}{15}$
t	0.30	0.49	0.69	0.91	1.17	1.47	
$P(t)$	$\frac{1}{15}$	$\frac{2}{15}$	$\frac{1}{15}$	$\frac{2}{15}$	$\frac{1}{15}$	$\frac{1}{15}$	

Permutation test, Rank test and Bootstrap

- Using a part of test statistics from the permutation distribution can get an estimate of p value (Bootstrap method) and sampling can be done with replacement. If M test statistics $t_i, i = 1, \dots, m$ are randomly sampled from the permutation distribution, the one-sided estimated p-value of Bootstrap is

$$\hat{p} = \frac{1 + \sum_{i=1}^M I(t_i \geq t)}{M + 1}$$

- Including the observed value t , there are $M+1$ test statistic values. Whereas the one-sided exact p-value of permutation test is

$$p = \frac{\sum_{i=1}^{\binom{N}{n}} I(t_i \geq t)}{\binom{N}{n}}$$

Permutation test, Rank test and Bootstrap

- p-value of Bootstrap is approximate
- p-value of permutation test is exact
- The rank test is a special permutation test, its test statistics does not depend on the observed data, but depend on the rank of the observed data, so it is more robust.

Multiple Testing

- The classical hypothesis testing problem is to consider limited parameters, say $\theta = \theta_0$ or no
- Sometimes, we need to do multiple testing at the same time, say

$$H_{10} : \theta_1 = 0 \quad \text{versus} \quad H_{11} : \theta_1 \neq 0$$

$$H_{20} : \theta_2 = 0 \quad \text{versus} \quad H_{21} : \theta_2 \neq 0$$

... ..

$$H_{N0} : \theta_N = 0 \quad \text{versus} \quad H_{N1} : \theta_N \neq 0$$

- For example, we want to identify the genes that cause one specific disease. For each gene, we want to know whether it works or not. The number of genes is pretty large here.
- Can we do the hypothesis testing problems individually, and reject H_0 if any individual H_0 is rejected?
- No and Yes.
- We want to control the overall false positives

Multiple Testing

- Recall. For a level α test, we have $\sup_{\theta \in \Theta_0} \beta(\theta) \leq \alpha$.
- Consider individual testing problem:

$$H_{i0} : \theta_i = 0 \quad \text{versus} \quad H_{i1} : \theta_i \neq 0$$

Suppose we have a level α test for this problem. Denote the power function of this test as $\beta_i(\theta)$. Then $\beta_i(\theta) \leq \alpha$.

- Therefore, for the original multiple testing problem, the rejection probability is

$$\beta(\theta) = 1 - P(\text{accept } H_{i0} \text{ for all } i \leq i \leq N) \stackrel{\text{indep}}{=} 1 - \prod_{i=1}^N (1 - \beta_i(\theta))$$

- Under null hypothesis $\theta = 0$, if all the test statistic are independent and all the tests have size α . The rejection probability is

$$1 - \prod_{i=1}^N (1 - \beta_i(0)) = 1 - (1 - \alpha)^N$$

How large it is? Let $N = 50, \alpha = 0.05$, then $1 - (1 - \alpha)^N \approx 0.92$.

$P(\text{rejection} | H_0) = 0.92!$

Familywise Error Control

- To control the overall false positives, we consider

$$P(\text{making at least one false rejection})$$

- Define $I = \{i; H_{i0} \text{ is true}\}$ be the index set for which H_0 is true
- Define $R = \{i; H_{j0} \text{ is rejected}\}$ be the index set that we reject.
- Define the familywise error rate at level α if

$$P(R \cap I \neq \emptyset) = P(\text{making at least one false rejection}) \leq \alpha.$$

- Bonferroni method:** for each individual test, set the level to be α/N .
Let p_j be the p -value for test H_{j0} versus H_{j1} .

$$\begin{aligned} P(\text{making a false rejection}) &= P(p_j < \alpha/N \text{ for some } i \in I) \\ &\leq \sum_{i \in I} P(p_j < \alpha/N) \\ &= \sum_{i \in I} \alpha/N = \frac{\alpha |I|}{N} \leq \alpha \end{aligned}$$

So we have overall control of the type I error.

- It can have low power.

Normal Example

Suppose we have N sample means Y_1, \dots, Y_N , each is the average of n normal observations with variance σ^2 . So $Y_j \sim N(\mu_j, \sigma^2/n)$. To test $H_{j0} : \mu_j = 0$ we can use the test statistic

$$T_j = \sqrt{n}Y_j/\sigma \sim N(\mu_j, 1).$$

The power function at $\mu_j = 0$ (p -value) is $p_j = 2\Phi(-|T_j|)$.

- If we did uncorrected testing that we reject $p_j < \alpha$, which means $|T_j| > z_{\alpha/2}$.
- With Bonferroni correction we reject when $p_j < \alpha/N$, which corresponds to

$$|T_j| > z_{\alpha/2N}$$

- Generate random samples under H_0 with code in next slide.
- If we apply the approximation for normal CDF and PDF, then

$$\frac{\phi(x)}{x + 1/x} \leq 1 - \Phi(x) \leq \frac{\phi(x)}{x}, \quad \phi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2},$$

then approximately, the corrected bound becomes $\sigma\sqrt{2\log(2N/\alpha)/n}$. It grows like $\sqrt{\log N}$.

R code

```

1 rm(list=ls())
2 N = 50; n = 100; sigma = 3; alpha = 0.05;
3 iter = 100; fwr1 = rep(0, iter); fwr2 = rep(0, iter);
4 for(i in 1:iter){
5   Y = rnorm(N, mean = 0, sd = sigma/sqrt(n));
6   #Generate Y_i's under null
7   stat = sqrt(n)*Y/sigma;
8   #Calculate the test statistic T_i
9   p = 1*(abs(stat) > qnorm(1 - alpha/2))
10  #Find the p-value for each individual test without correction
11  corp = 1*(abs(stat) > qnorm(1 - alpha/2/N))
12  #Find the p-value for each individual test with Bonferroni co
13  fwr1[i] = 1*(sum(p) > 0); #Familywise error for test 1;
14  fwr2[i] = 1*(sum(corp) > 0); #Familywise error for test 2;
15 }
16 mean(fwr1) #empirical familywise error for uncorrected test #0.
17 mean(fwr2) #empirical familywise error for corrected test #0.06
18
19

```


Higher Criticism

- Let p_j be the p -value for test problem H_j . Then under null hypothesis, $p_j \sim \text{Unif}(0, 1)$.
- Consider a level α test for an individual hypothesis test, we reject the hypothesis when $p_j < \alpha$
- Let $Y_j = I(p_j \leq \alpha)$. Under null hypothesis, $Y_j \sim \text{Bernoulli}(\alpha)$. So \bar{Y}_n has mean α and standard deviation $\sqrt{\alpha(1-\alpha)/N}$
- According to CLT, let \bar{Y}_N be the fraction of the rejected hypothesis with a level α test, then

$$T_N = \sqrt{N} \frac{\bar{Y}_n - \alpha}{\sqrt{\alpha(1-\alpha)}} \xrightarrow{d} N(0, 1),$$

When $|T_n| > z_{\alpha/2}$, at least one hypothesis testing problem rejects H_0

- Why do we consider a fixed α ? When α changes, will we get different results?
- Define the function of α . The expression is as following:
 - (1) Sort p -values so that $p_{(1)} < p_{(2)} < \dots < p_{(N)}$
 - (2) Define

$$HC_k = \sqrt{N} \frac{k/N - p_{(k)}}{\sqrt{p_{(k)}(1 - p_{(k)})}}.$$

- (3) The test statistic is $T_N = \max_{1 \leq k \leq N} HC_k$

 - Note. Now $p_{(k)}$ plays the role of α . If we take $\alpha = p_{(k)}$, then the number of rejected hypothesis is k , so the fraction is k/N .
 - The limiting distribution for T_n is Gumbel distribution (no need to know)

