

# Financial Time Series I

## Topic 3: Resampling Methods

Hung Chen

Department of Mathematics

National Taiwan University

11/11/2002

# OUTLINE

1. Motivated Example
  - Double Blind Randomized Experiment
  - Bootstrap with R
2. Odds Ratio
  - Definition
  - Random Design
  - Prospective Study
  - Retropective Study
3. Bootstrap Method
  - Parametric Bootstrap
  - Nonparametric Bootstrap
  - Failure of bootstrap method

## The Practice of Statistics

- Statistics is the science of learning from experience, especially experience that arrives a little bit at a time.
- Most people are not natural-born statisticians.
  - We are not very good at picking out patterns from a sea of noisy data.
  - To put it another way, we all are too good at picking out non-existent patterns that happen to suit our purposes.
  - Statistical theory attacks the problem from both ends. It provides optimal methods for finding a real signal in a noisy background, and also provides strict checks against the overinterpretation of random patterns.
- Statistical theory attempts to answer three basic questions:
  1. Data Collection: How should I collect my data?
  2. Summary: How should I analyze and summarize the data that I've collected?

### 3. Statistical Inference: How accurate are my data summaries?

- The bootstrap is a recently developed technique for making certain kinds of statistical inferences.

It is only recently developed because it requires modern computer power to simplify the often intricate calculations of traditional statistical theory.

- The idea of bootstrap method is close to that of simulation. The main difference is to plug in an estimate of the underlying unknown random mechanism  $F$ .

## Motivated Example

- Illustrate the just mentioned three basic statistical concepts using a front-page news from the New York Times of January 27, 1987.
- A study was done to see if small aspirin doses would prevent heart attacks in healthy middle-aged men.
- The data for the aspirin study were collected in a particularly efficient way: by a controlled, randomized, double-blind study.
  - One half of the subjects received aspirin and the other half received a control substance, or placebo, with no active ingredients.
  - The subjects were randomly assigned to the aspirin or placebo groups.
  - Both the subjects and the supervising physicians were blind to the assignments, with the statisticians keeping a secret code of who received which substance.
  - Scientists, like everyone else, want the subject they are working on to succeed.

- The elaborate precautions of a controlled, randomized, blinded experiment guard against seeing benefits that don't exist, while maximizing the chance of detecting a genuine positive effect.
- The summary statistics in the study are very simple:

	heart attacks (fatal plus non-fatal)	subjects
aspirin group:	104	11,037
placebo group:	189	11,034

- What strikes the eye here is the lower rate of heart attacks in the aspirin group.
- The ratio of the two rates is

$$\hat{\theta} = \frac{104/11037}{189/11034} = 0.55.$$

It suggests that the aspirin-takers only have 55% as many as heart attacks as placebo-takers.

- We are not interested in  $\hat{\theta}$ . What we would like to know is  $\theta$ , the true ratio, that is the ratio we would see if we

could treat all subjects, and not just a sample of them.

- The tough question is how do we know that  $\hat{\theta}$  might not come out much less favorably if the experiment were run again? This is where statistical inference comes in.
- Statistical theory allows us to make the following inference: the true value of  $\theta$  lies in the interval  $0.43 < \theta < 0.70$  with 95% confidence.

Note that

$$\theta = \hat{\theta} + (\theta - \hat{\theta}) = 0.55 + [\theta - \hat{\theta}(\omega_0)],$$

where  $\theta$  and  $\hat{\theta}(\omega_0)$  ( $= 0.55$ ) are two numbers.

- Since  $\omega_0$  cannot be observed, we use  $\theta - \hat{\theta}(\omega)$  to describe  $\theta - \hat{\theta}(\omega_0)$  in statistics.
- What is the fluctuation of  $\theta - \hat{\theta}(\omega)$  among all  $\omega$ ?
- If, for most  $\omega$ ,  $\theta - \hat{\theta}(\omega)$  is around zero, we can conclude statistically that  $\theta$  is close to 0.55 ( $= \hat{\theta}(\omega_0)$ ).

- (Recall the definition of consistency.) If

$$P(\omega : |\theta - \hat{\theta}(\omega)| < 0.1) = 0.95,$$

we claim that with 95% confidence that  $\theta - 0.55$  is no more than 0.1.

- In the aspirin study, it also track strokes. The results are presented as the following:

	strokes	subjects
aspirin group:	119	11,037
placebo group:	98	11,034

- For strokes, the ratio of the two rates is

$$\hat{\theta} = \frac{119/11037}{98/11034} = 1.21.$$

It now looks like taking aspirin is actually harmful.

- However, the interval for the true stroke ratio  $\theta$  turns out to be  $0.93 < \theta < 1.59$  with 95% confidence. This includes the neutral value  $\theta = 1$ , at which aspirin would be no better or worse than placebo.
- In the language of statistical hypothesis testing, aspirin was found to be significantly



beneficial for preventing heart attacks, but not significantly harmful for causing strokes.

According to introductory statistics course, we can put the above problem into the framework of two-sample problem with binomial distribution. Asymptotic analysis will then be used to give an approximation on the distribution of  $\hat{\theta}$ .

In this note, we demonstrate an alternative. Apply the bootstrap method in the stroke example.

1. Create two *pseudo* populations based on the collected data:
  - Pseudo population 1: It is consist of 119 ones and  $11037 - 119 = 10918$  zeros.
  - Pseudo population 2: It is consist of 98 ones and  $11034 - 98 = 10936$  zeros.
2. (Monte Carlo Resampling) Draw with replacement a sample of 11037 items from the first pseudo population, and a sample of 11034 items from the second pseudo population. Each of these is called a

*bootstrap sample.*

3. Derive the bootstrap replicate of  $\hat{\theta}$ :

$$\hat{\theta}^* = \frac{\text{prop. of ones in bootstrap sample \#1}}{\text{prop. of ones in bootstrap sample \#2}}$$

4. Repeat this process (1-3) a large number of times, say 1000 times, and obtain 1000 *bootstrap replicates*  $\hat{\theta}^*$ .

The above procedure can be implemented easily using the following R code.

- $n1 < -11037$ ;  $s1 < -119$ ;  $p1 < -s1/n1$
- $n2 < -11034$ ;  $s2 < -98$ ;  $p2 < -s2/n2$
- Write a function named `stroke`.
- ```
stroke <- function(n1, p1, n2, p2){  
  control <- rbinom(1, n1, p1)  
  treat <- rbinom(1, n2, p2)  
  theta <- -(control/n1)/(treat/n2)  
  return(theta)  
}
```
- Suppose that we would like to do 1000 replications.

- $result < -rep(1, 1000)$  which is used to store the 1000 bootstrap replicates of  $\hat{\theta}$ .
- for  $(i \text{ in } 1 : 1000)$   $result[i] < -stroke(n1, p1, n2, p2)$

My simulation gives

- The standard deviation turned out to be 0.17 in a batch of 1000 replicates that we generated.
- A rough 95% confidence interval is (0.93, 1.60) which is derived by taking the 25th and 975th largest of the 1000 replicates.
- The above method is called the percentile method.

**Project 1.1** Do a computer experiment to repeat the above process one hundred times and answer the following questions.

- (a) Give a summary of those 100 confidence intervals. Do they always contain 1? If not, is it *statistical correct*?
- (b) How do you describe the distribution of bootstrap replicates of  $\hat{\theta}$ ? Is it close to normal?

## Odds Ratio

- If an event has probability  $P(A)$  of occurring, the **odds** of  $A$  occurring are defined to be

$$\text{odds}(A) = \frac{P(A)}{1 - P(A)}.$$

- Let  $X$  denote the event that an individual is exposed to a potentially harmful agent and  $D$  denote the event that the individual becomes diseased.

Denote the complementary events as  $\bar{X}$  and  $\bar{D}$ .

- The odds of an individual contracting the disease given that he is exposed are

$$\text{odds}(D|X) = \frac{P(D|X)}{1 - P(D|X)}$$

and the odds of contracting the disease given that he is not exposed are

$$\text{odds}(D|\bar{X}) = \frac{P(D|\bar{X})}{1 - P(D|\bar{X})}.$$

- The **odds ratio**  $\Delta = \frac{\text{odds}(D|X)}{\text{odds}(D|\bar{X})}$  is a measure of the influence of exposure on subsequent disease.

We will consider how the odds and odds ratio could be estimated by sampling from a population with joint and marginal probabilities defined as in the following table:

|           |            |            |            |
|-----------|------------|------------|------------|
|           | $\bar{D}$  | $D$        |            |
| $\bar{X}$ | $\pi_{00}$ | $\pi_{01}$ | $\pi_{0.}$ |
| $X$       | $\pi_{10}$ | $\pi_{11}$ | $\pi_{1.}$ |
|           | $\pi_{.0}$ | $\pi_{.1}$ | 1          |

With this notation,

$$P(D|X) = \frac{\pi_{11}}{\pi_{10} + \pi_{11}} \quad P(D|\bar{X}) = \frac{\pi_{01}}{\pi_{00} + \pi_{01}}$$

so that

$$odds(D|X) = \frac{\pi_{11}}{\pi_{10}} \quad odds(D|\bar{X}) = \frac{\pi_{01}}{\pi_{00}}$$

and the odds ratio is

$$\Delta = \frac{\pi_{11}\pi_{00}}{\pi_{01}\pi_{10}}$$

the product of the diagonal probabilities in the preceding table divided by the product of the off-diagonal probabilities.

Now we will consider three possible ways to sample this population to study the relationship of disease and exposure.

- Random sample:
  - From such a sample, we could estimate all the probabilities directly.
  - If the disease is rare, the total sample size would have to be quite large to guarantee that a substantial number of diseased individuals was included.
- Prospective study:
  - A fixed number of exposed and nonexposed individuals are sampled and then followed through time.
  - The incidences of disease in those two groups are compared.
  - In this case the data allow us to estimate and compare  $P(D|X)$  and  $P(D|\bar{X})$  and, hence, the odds ratio.
  - The aspirin study described in the previous section can be viewed as this type of study.
- Retrospective study:
  - A fixed number of diseased and undiseased individuals are sampled and the

incidences of exposure in the two groups are compared.

- From such data we can directly estimate  $P(X|D)$  and  $P(X|\bar{D})$ .
- Since the marginal counts of diseased and nondiseased are fixed, we cannot estimate the joint probabilities or the important conditional probabilities  $P(D|X)$  and  $P(D|\bar{X})$ .
- Observe that

$$\begin{aligned}
 P(X|D) &= \frac{\pi_{11}}{\pi_{01} + \pi_{11}}, \\
 1 - P(X|D) &= \frac{\pi_{01}}{\pi_{01} + \pi_{11}}, \\
 odds(X|D) &= \frac{\pi_{11}}{\pi_{01}}, \\
 odds(X|\bar{D}) &= \frac{\pi_{10}}{\pi_{00}}.
 \end{aligned}$$

The odds ratio can also be expressed as  $odds(X|D)/odds(X|\bar{D})$ .

Now we describe the study of Vianna, Greenwald, and Davies (1971) to illustrate the retrospective study.

- In this study they collected data comparing the percentages of tonsillectomies for a

group of patients suffering from Hodgkin's disease and a comparable control group:

|           | Tonsillectomy | No Tonsillectomy |
|-----------|---------------|------------------|
| Hodgkin's | 67            | 34               |
| Control   | 43            | 64               |

- Recall that the odds ratio can be expressed as  $odds(X|D)/odds(X|\bar{D})$  and an estimate of it is  $n_{00}n_{11}/(n_{01}n_{10})$ , the product of the diagonal counts divided by the product of the off-diagonal counts.
- The data of Vianna, Greenwald, and Davies gives an estimate of odds ratio is

$$\frac{67 \times 64}{43 \times 34} = 2.93.$$

- According to this study, the odds of contracting Hodgkin's disease is increased by about a factor of three by undergoing a tonsillectomy.
- As well as having a point estimate 2.93, it would be useful to attach an approximate standard error to the estimate to indicate its uncertainty.



- We will use simulation (parametric bootstrap) to approximate the distribution of  $\Delta$ .
  - We need to generate random numbers according to a statistical model for the counts in the table of Vianna, Greenwald, and Davies.
  - The model is that the count in the first row and first column,  $N_{11}$ , is binomially distributed with  $n = 101$  and probability  $\pi_{11}$ .
  - The count in the second row and second column,  $N_{22}$ , is binomially distributed with  $n = 107$  and probability  $\pi_{22}$ .
  - The distribution of the random variable

$$\hat{\Delta} = \frac{N_{11}N_{22}}{(101 - N_{11})(107 - N_{22})}$$

is thus determined by the two binomial distributions, and we could approximate it arbitrarily well by drawing a large number of samples from them.

- Since the probabilities  $\pi_{11}$  and  $\pi_{22}$  are unknown, they are estimated from the observed counts by  $\hat{\pi}_{11} = 67/101 = 0.663$

and  $\pi_{22} = 64/107 = 0.598$ .

- A one thousand realizations generated on a computer gives the standard deviation 0.89.

**Project 1.2** Do a computer experiment to run the above process in the setting of retrospective study. Give a 95% confidence interval of  $\Delta$  and describe the distribution of bootstrap replicates of  $\hat{\Delta}$ ?

## Bootstrap Method

- The bootstrap method introduced in Efron (1979) is a very general resampling procedure for estimating the distributions of statistics based on independent observations.
  - The bootstrap method is shown to be successful in many situations, which is being accepted as an alternative to the asymptotic methods.
  - It is better than some other asymptotic methods, such as the traditional normal approximation and the Edgeworth expansion.
  - There are some counterexamples that show the bootstrap produces wrong solutions, i.e., it provides some inconsistent estimators.

Consider the problem of estimating variability of location estimates by the Bootstrap method.

- If we view the observations  $x_1, x_2, \dots, x_n$  as realizations of independent random variables with common distribution function  $F$ ,

it is appropriate to investigate the variability and sampling distribution of a location estimate calculated from a sample of size  $n$ .

- Denote the location estimate as  $\hat{\theta}$ .
  - Note that  $\hat{\theta}$  is a function of the random variables  $X_1, X_2, \dots, X_n$  and hence has a probability distribution, its sampling distribution, which is determined by  $n$  and  $F$ .
  - How do we derive this sampling distribution?
- We are faced with two problems:
  1.  $F$  is unknown.
  2.  $F$  is known, but  $\hat{\theta}$  may be such a complicated function of  $X_1, X_2, \dots, X_n$  that finding its distribution would exceed our analytic abilities.
- To address the second problem when  $F$  is known.
  - How could we find the probability distribution of  $\hat{\theta}$  without going through incredibly complicated analytic calculations?

- The computer comes to our rescue—we can do it by simulation.
- We generate many, many samples, say  $B$  in number, of size  $n$  from  $F$ ; from each sample we calculate the value of  $\hat{\theta}$ .
- The empirical distribution of the resulting values  $\hat{\theta}_1^*, \hat{\theta}_2^*, \dots, \hat{\theta}_B^*$  is an approximation to the distribution function of  $\hat{\theta}$ , which is good if  $B$  is very large.
- If we wish to know the standard deviation of  $\hat{\theta}$ , we can find a good approximation to it by calculating the standard deviation of the collection of values  $\hat{\theta}_1^*, \hat{\theta}_2^*, \dots, \hat{\theta}_B^*$ .
- We can make these approximations arbitrarily accurate by taking  $B$  to be arbitrarily large.

**Simulation** Let  $G$  be a distribution and let  $Y_1, \dots, Y_B$  be iid values drawn from  $G$ .

- By the law of large numbers,  $B^{-1} \sum_{j=1}^B Y_j$  converges in probability to  $E(Y)$ .
- We can use  $B^{-1} \sum_{j=1}^B Y_j$  as an estimate of  $E(Y)$ .

- In a simulation, we can make  $B$  as large as we like in which case, the difference between  $B^{-1} \sum_{j=1}^B Y_j$  and  $E(Y)$  is negligible.

All this would be well and good if we knew  $F$ , but we don't. So what do we do? We will consider two different cases.

- In the first case,  $F$  is unknown up to an unknown parameter  $\eta$ , i.e.  $F(x|\eta)$ .
  - Without knowing  $\eta$ , the above approximation cannot be used.
  - The idea of the **parametric bootstrap** is to simulate data from  $F(x|\hat{\eta})$  where  $\hat{\eta}$  should be a good estimate of  $\eta$ .
  - It utilizes the structure of  $F$ .
- In the second case,  $F$  is completely unknown.
- The idea of the **nonparametric bootstrap** is to simulate data from the empirical cdf  $F_n$ .
- Here  $F_n$  is a discrete probability distribution that gives probability  $1/n$  to each observed value  $x_1, \dots, x_n$ .

- A sample of size  $n$  from  $F_n$  is thus a sample of size  $n$  drawn *with replacement* from the collection  $x_1, \dots, x_n$ . The standard deviation of  $\hat{\theta}$  is then estimated by

$$s_{\hat{\theta}} = \sqrt{\frac{1}{B} \sum_{i=1}^B (\theta_i^* - \bar{\theta}^*)^2}$$

where  $\theta_1^*, \dots, \theta_B^*$  are produced from  $B$  sample of size  $n$  from the collection  $x_1, \dots, x_n$ .

Now we use a simple example to illustrate this idea.

- Suppose  $n = 2$  and observe  $X_{(1)} = c < X_{(2)} = d$ .
- $X_1^*, X_2^*$  are independently distributed with  $P(X_i^* = c) = P(X_i^* = d) = 1/2, \quad i = 1, 2$ .
- The pairs  $(X_1^*, X_2^*)$  therefore takes on the four possible pairs of values

$$(c, c), (c, d), (d, c), (d, d),$$

each with probability  $1/4$ .

- $\theta^* = (X_1^* + X_2^*)/2$  takes on the values  $c, (c+d)/2, d$  with probabilities  $1/4, 1/2, 1/4$ ,

respectively, so that  $\theta^* - (c + d)/2$  takes on the values  $(c - d)/2$ ,  $0$ ,  $(d - c)/2$  with probabilities  $1/4$ ,  $1/2$ ,  $1/4$ , respectively.

For the above example, we can easily calculate its bootstrap distribution.

- When  $n$  is large, we can easily imagine that the above computation becomes too complicated to compute directly.
- Use simple random sampling to approximate bootstrap distribution.
- In the bootstrap literature, a variety alternatives are suggested other than simple random sampling.

**Project 1.3** Use parametric bootstrap and nonparametric bootstrap to approximate the distribution of median based on a data with sample size 20 from a standard normal distribution. The following is a sample R-code.

- $n < -20$
- $x < -rnorm(n)$  # Create some data.
- $theta.hat < -median(x)$



- $B < -1000$ ; theta.boot  $< -$  rep(0,B)
- for( $i$  in 1 :  $B$ ) {  
xstar  $< -$  sample( $x$ ,size= $n$ ,replace=T) #  
draw a bootstrap sample  
theta.boot[ $i$ ]  $< -$  median(xstar) # com-  
pute the statistic  
}
- var.boot  $< -$  var(theta.boot)
- se|- sqrt(var.boot); print(se)

We now introduce notations to illustrate the bootstrap method.

- Assumed the data  $X_1, \dots, X_n$ , are independent and identically distributed (iid) samples from a  $k$ -dimensional population distribution  $F$ .
- Estimate the distribution

$$H_n(x) = P\{R_n \leq x\},$$

where  $R_n = R_n(T_n, F)$  is a real-valued functional of  $F$  and  $T_n = T_n(X_1, \dots, X_n)$ , a statistic of interest.

- Let  $X_1^*, \dots, X_n^*$  be a “bootstrap” samples iid from  $F_n$ , the empirical distribution based on  $X_1, \dots, X_n$ ,  $T_n^* = T_n(X_1^*, \dots, X_n^*)$ , and  $R_n^* = R_n(T_n^*, F_n)$ .  $F_n$  is constructed by placing at each observation  $X_i$  a mass  $1/n$ . Thus  $F_n$  may be represented as

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x), \quad -\infty < x < \infty.$$

- A bootstrap estimator of  $H_n$  is

$$\hat{H}_n(x) = P_*\{R_n^* \leq x\},$$

where for given  $X_1, \dots, X_n$ ,  $P_*$  is the conditional probability with respect to the random generator of bootstrap samples.

- Since the bootstrap samples are generated from  $F_n$ , this method is called the nonparametric bootstrap.
  - Note that  $\hat{H}_n(x)$  will depend on  $F_n$  and hence itself is a random variable.
  - To be specific,  $\hat{H}_n(x)$  will change as the data  $\{x_1, \dots, x_n\}$  changes.
  - Recall that a bootstrap analysis is run to assess the accuracy of some primary statistical results.
  - This produces bootstrap statistics, like standard errors or confidence intervals, which are assessments of error for the primary results.
- As a further remark, the empirical distribution  $F_n$  is called the nonparametric maximum likelihood estimate (MLE) of  $F$ .

As illustration, we consider the following three examples.

**Example 1.** Suppose that  $X_1, \dots, X_n \sim N(\mu, 1)$  and  $R_n = \sqrt{n}(\bar{X}_n - \mu)$ . Consider the estimation of

$$P(a) = P\{R_n > a | N(\mu, 1)\}.$$

The nonparametric bootstrap method will estimate  $P(a)$  by

$$P_{NB}(a) = P\{\sqrt{n}(\bar{X}_n^* - \bar{X}_n) > a | F_n\}.$$

- Observe data  $x_1, \dots, x_n$  with mean  $\bar{x}_n$ .
- Let  $Y_1, \dots, Y_n$  denote a bootstrap sample of  $n$  observations drawn independently from  $F_n$ .
- Let  $\bar{Y}_n = n^{-1} \sum_{i=1}^n Y_i$ .
- $P(a)$  is estimated by

$$P_{NB}(a) = P\{\sqrt{n}(\bar{Y}_n - \bar{x}_n) > a | F_n\}.$$

- In principle,  $P_{NB}(a)$  can be found by considering all  $n^n$  possible bootstrap sample.
  - If all  $X_i$ 's are distinct, then the number of different possible resamples equals the number of distinct ways of placing  $n$  indistinguishable objects into  $n$  numbered

- boxes, the boxes being allowed to contain any number of objects. It is known that it is equal to  $C(2n-1, n) \approx (n\pi)^{-1/2}2^{2n-1}$ .
- When  $n = 10(20, \text{ respect.})$ ,  $C(2n-1, n) \approx 92375(6.9 \times 10^{10}, \text{ respect.})$ .
  - For small value of  $n$ , it is often feasible to calculate a bootstrap estimate exactly.
  - For large samples, say  $n \geq 10$ , this becomes infeasible even at today’s computer technology.
- Natural questions to ask are as follows:
    - What are computationally efficient ways to bootstrap?
    - Can we get bootstrap-like answers without Monte Carlo?
  - Address the question of “evaluating” the performance of bootstrap method.
    - For the above particular problem, we need to estimate  $P_{NB}(a) - P(a)$  or  $\sup_a |P_{NB}(a) - P(a)|$ .
    - As a remark,  $P_{NB}(a)$  is a random variable since  $F_n$  is random.

- Efron (1992) proposed to use jackknife to give the error estimates for bootstrap quantities.
- Suppose that additional information on  $F$  is available. Then it is reasonable to utilize this information in the bootstrap method.
- In this example,  $F$  known to be normally distributed with unknown mean  $\mu$  and variance 1.
  - It is natural to use  $\bar{x}_n$  to estimate  $\mu$  and then estimate  $P(a) = P\{R_n > a | N(\mu, 1)\}$  by
 
$$P_{PB}(a) = P\{\sqrt{n}(\bar{Y}_n - \bar{x}_n) > a | N(\bar{x}_n, 1)\}.$$
  - Since the bootstrap samples are generated from  $N(\bar{x}_n, 1)$  which utilizes the information from a parametric form of  $F$ , this method is called the parametric bootstrap.
  - In this case, it can be shown that  $P_{PB}(a) = P(a)$  for all realization of  $\bar{X}_n$ .
  - If  $F$  is known to be normally distributed with unknown mean and variance  $\mu$  and

variance  $\sigma^2$  respectively,  $P_{PB}(a)$  is no longer equal to  $P(a)$ .

**Project 1.4.** (a) Show that  $P_{PB}(a) = \Phi(a/s_n)$  where  $s_n^2 = (n - 1)^{-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2$ .

(b) Prove that  $P_{PB}(a)$  is a consistent estimate of  $P(a)$  for fixed  $a$ .

(c) Prove that  $\sup_a |P_{PB}(a) - P(a)| \xrightarrow{P} 0$ .

For the question of finding  $P_{NB}(a)$ , we can in principle write down the characteristic function and then apply the inversion formula. However, it is a nontrivial job. Therefore, Efron (1979) suggested to approximate  $P_{NB}(a)$  by Monte Carlo resampling. (i.e., Sample-size resamples may be drawn repeatedly from the original sample, the value of a statistic computed for each individual resample, and the bootstrap statistic approximated by taking an average of an appropriate function of these numbers.)

Now we state Levy's Inversion Formula which is taken from Chapter 6.2 of Chung (1974).

**Theorem** If  $x_1 < x_2$  and  $x_1$  and  $x_2$  are points of continuity of  $F$ , then we have

$$F(x_2) - F(x_1) = \lim_{T \rightarrow \infty} \frac{1}{2\pi} \int_{-T}^T \frac{e^{-itx_1} - e^{-itx_2}}{it} f(t) dt,$$

where  $f(t)$  is the characteristic function.



Example 2. Estimating the probability of  
success

- Consider a probability distribution  $F$  putting all of its mass at zero or one.
- Let  $\theta(F) = P(X = 1) = p$ .
- Consider  $R(\mathbf{X}, F) = \bar{X} - \theta(F) = \hat{p} - p$ .
- Observed  $\mathbf{X} = \mathbf{x}$ , the bootstrap sample  $X_1^*, \dots, X_n^* \sim \text{Bin}(1, \theta(F_n)) = \text{Bin}(1, \bar{x}_n)$ .

Note that

$$\begin{aligned} R(\mathbf{X}^*, F_n) &= \bar{X}_n^* - \bar{x}_n, \\ E_*(\bar{X}_n^* - \bar{x}_n) &= 0, \\ \text{Var}_*(\bar{X}_n^* - \bar{x}_n) &= \frac{\bar{x}_n(1 - \bar{x}_n)}{n}. \end{aligned}$$

Recall that  $n\bar{X}_n^* \sim \text{Bin}(n, \bar{x})$  and  $n\bar{X}_n \sim \text{Bin}(n, p)$ .

- It is known that if  $\min\{n\bar{x}_n, n(1 - \bar{x}_n)\} \geq 5$ ,

$$\frac{n\bar{X}_n^* - n\bar{x}_n}{\sqrt{n\bar{x}_n(1 - \bar{x}_n)}} = \frac{\sqrt{n}(\bar{X}_n^* - \bar{x}_n)}{\sqrt{\bar{x}_n(1 - \bar{x}_n)}} \sim N(0, 1);$$

and if  $\min\{np, n(1 - p)\} \geq 5$ ,

$$\frac{n\bar{X}_n - np}{\sqrt{n\theta(1 - p)}} = \frac{\sqrt{n}(\bar{X}_n - p)}{\sqrt{p(1 - p)}} \sim N(0, 1).$$

- Based on the above approximation results, we conclude that the bootstrap method works if  $\min\{n\bar{x}_n, n(1 - \bar{x}_n)\} \geq 5$ .
- The question remained to be studied is whether

$$P\{\min(n\bar{X}_n, n(1 - \bar{X}_n)) \geq 5\} \rightarrow 0?$$

### Example 3. Estimating the median

- Suppose we are interested in finding the distribution of  $n^{1/2}\{F_n^{-1}(1/2) - F^{-1}(1/2)\}$  where  $F_n^{-1}(1/2)$  and  $F^{-1}(1/2)$  are the sample and population median respectively.
- Set  $\theta(F) = F^{-1}(1/2)$ .
- Find a bootstrap approximation of the above distribution.
- Consider  $n = 2m - 1$ . Then the sample median  $F_n^{-1}(1/2) = X_{(m)}$  where  $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ .
- Let  $N_i^*$  denote the number of times  $x_i$  is selected in the bootstrap sampling procedure. Set  $\mathbf{N}^* = (N_1^*, \dots, N_n^*)$ . It follows easily that  $\mathbf{N}^*$  follows a multinomial distribution with  $n$  trials and the probability of selection is  $(n^{-1}, \dots, n^{-1})$ .
- Denote the order statistics of  $x_1, \dots, x_n$  by  $x_{(1)} \leq \dots \leq x_{(n)}$ .
- Set  $N_{[i]}^*$  to be the number of times of choosing  $x_{(i)}$ . Then for  $1 \leq \ell < n$ , we have
 
$$Prob_*(X_{(m)}^* > x_{(\ell)}) = Prob_*\{N_{[1]}^* + \dots + N_{[\ell]}^* \leq m - 1\}$$

$$\begin{aligned}
&= \text{Prob} \left\{ \text{Bin} \left( n, \frac{\ell}{n} \right) \leq m - 1 \right\} \\
&= \sum_{j=0}^{m-1} C(n, j) \left( \frac{\ell}{n} \right)^j \left( 1 - \frac{\ell}{n} \right)^{n-j}.
\end{aligned}$$

Or,

$$\begin{aligned}
\text{Prob}_*(T^* = x_{(\ell)} - x_{(m)}) &= \text{Prob} \left\{ \text{Bin} \left( n, \frac{\ell - 1}{n} \right) \leq m - 1 \right. \\
&\quad \left. - \text{Prob} \left\{ \text{Bin} \left( n, \frac{\ell}{n} \right) \leq m - 1 \right\} \right\}.
\end{aligned}$$

- When  $n = 13$ , we have

|             |         |         |         |        |        |       |
|-------------|---------|---------|---------|--------|--------|-------|
| $\ell$      | 2 or 12 | 3 or 11 | 4 or 10 | 5 or 9 | 6 or 8 | 7     |
| probability | 0.0015  | 0.0142  | 0.0550  | 0.1242 | 0.4136 | 0.223 |

Quite often we use the mean square error to measure the performance of an estimator,  $t(X)$ , of  $\theta(F)$ . Or,  $E_F T^2 = E_F (t(X) - \theta(F))^2$ . Use the bootstrap to estimate  $E_F T^2$ . Then the bootstrap estimate of  $E_F T^2$  is

$$E_*(T^*)^2 = \sum_{\ell=1}^{13} [x_{(\ell)} - x_{(7)}]^2 \text{Prob}_* \{T^* = x_{(\ell)} - x_{(7)}\}.$$

It is known that  $E_F T^2 \rightarrow [4nf^2(\theta)]^{-1}$  as  $n$  tends to infinity when  $F$  has a bounded continuous density. A natural question to ask is whether  $E_*(T^*)^2$  is close to  $E_F T^2$ ?

## Validity of the Bootstrap Method

We now give a brief discussion on the validity of the bootstrap method. First, we state central limit theorems and its approximation error bound.

Perhaps the most widely known version of the CLT is the following.

**Theorem** (Lindeberg-Levy Central Limit Theorem) Let  $\{X_i\}$  be iid with mean  $\mu$  and finite variance  $\sigma^2$ . Then

$$\sqrt{n} \left( \frac{1}{n} \sum_{i=1}^n X_i - \mu \right) \xrightarrow{d} N(0, \sigma^2).$$

The above theorem can be generalized to independent random variables which are not necessarily identically distributed.

**Theorem** (Lindeberg-Feller CLT) Let  $\{X_i\}$  be independent with mean  $\{\mu_i\}$ , finite variances  $\{\sigma_i^2\}$ , and distribution functions  $\{F_i\}$ .

- Suppose that  $B_n^2 = \sum_{i=1}^n \sigma_i^2$  satisfies

$$\frac{\sigma_n^2}{B_n^2} \rightarrow 0, \quad B_n \rightarrow \infty \quad \text{as } n \rightarrow \infty.$$

- $n^{-1} \sum_{i=1}^n X_i$  is  $N(n^{-1} \sum_{i=1}^n \mu_i, n^{-2} B_n^2)$  if and

only if the following Lindeberg condition satisfied

$$B_n^{-2} \sum_{i=1}^n \int_{|t-\mu_i| > \epsilon B_n} (t-\mu_i)^2 dF_i(t) \rightarrow 0, \quad n \rightarrow \infty, \quad \text{each } \epsilon > 0$$

In the theorems just described, asymptotic normality was asserted for a sequence of sums  $\sum_1^n X_i$  generated by a single sequence  $X_1, X_2, \dots$  of random variables. For the validity of bootstrap, we may consider a *double array* of random variables

$$\begin{array}{cccc} X_{11}, & X_{12}, & \cdots, & X_{1K_1}; \\ X_{21}, & X_{22}, & \cdots, & X_{2K_2}; \\ \vdots & \vdots & \vdots & \vdots \\ X_{n1}, & X_{n2}, & \cdots, & X_{nK_n}; \\ \vdots & \vdots & \vdots & \vdots \end{array}$$

For each  $n \geq 1$ , there are  $K_n$  random variables  $\{X_{nj}, 1 \leq j \leq K_n\}$ . It is assumed that  $K_n \rightarrow \infty$ . The case  $K_n = n$  is called a “triangular” array.

Denote by  $F_{nj}$  the distribution function of  $X_{nj}$ . Also, put

$$\begin{aligned} \mu_{nj} &= EX_{nj}, \\ A_n &= E \sum_{j=1}^{K_n} X_{nj} = \sum_{j=1}^{K_n} \mu_{nj}, \end{aligned}$$

$$B_n^2 = \text{Var} \left( \sum_{j=1}^{K_n} X_{nj} \right).$$

We then have the following theorem.

**Theorem** (Lindeberg-Feller) Let  $\{X_{nj} : 1 \leq j \leq K_n; n = 1, 2, \dots\}$  be a double array with independent random variables within rows. Then the “uniform asymptotic negligibility” condition

$$\max_{1 \leq j \leq K_n} P(|X_{nj} - \mu_{nj}| > \tau B_n) \rightarrow 0, \quad n \rightarrow \infty, \text{ each } \tau > 0,$$

and the asymptotic normality condition  $\sum_{j=1}^{K_n} X_{nj}$  is  $AN(A_n, B_n^2)$  together hold if and only if the Lindberg condition

$$B_n^{-2} \sum_{i=1}^n \int_{|t - \mu_i| > \epsilon B_n} (t - \mu_i)^2 dF_i(t) \rightarrow 0, \quad n \rightarrow \infty \text{ each } \epsilon > 0$$

is satisfied. As a note, the independence is assumed only within rows, which themselves may be arbitrarily dependent.

It is of both theoretical and practical interest to characterize the error of approximation in the CLT.

For the i.i.d. case, an exact bound on the error of approximation is provided by the following theorem due to Berry (1941) and Esseen (1945).

**Theorem** If  $X_1, \dots, X_n$  are i.i.d. with distribution  $F$  and if  $S_n = X_1 + \dots + X_n$ , then there exists a constant  $c$  (independent of  $F$ ) such that for all  $x$ ,

$$\sup_x \left| P \left[ \frac{S_n - ES_n}{\sqrt{Var(S_n)}} \leq x \right] - \Phi(x) \right| \leq \frac{c}{\sqrt{n}} \frac{E|X_1 - EX_1|^3}{[Var(X_1)]^{3/2}}$$

for all  $F$  with finite third moment.

- Note that  $c$  in the above theorem is a universal constant. Various authors have thought to find the best constant  $c$ .
- Originally,  $c$  is set to be  $33/4$  but it has been sharpened to  $0.7975$ .
- For  $x$  is sufficiently large, while  $n$  remains fixed, the quantities

$$P[(S_n - ES_n)/\sqrt{Var(S_n)} \leq x]$$

and  $\Phi(x)$  each become so close to 1 that the bound given by above is too crude.

- The problem in this case may be characterized as one of approximation of *large deviation* probabilities, with the object of attention becoming the relative error in approxi-



mation of

$$1 - P[(S_n - ES_n)/\sqrt{Var(S_n)} \leq x]$$

by  $1 - \Phi(x)$  when  $x \rightarrow \infty$ .

## Inconsistent Bootstrap Estimator

Bickel and Freedman (1981) and Loh (1984) showed that the bootstrap estimators of the distributions of the extreme-order statistics are inconsistent.

- Let  $X_{(n)}$  be the maximum of i.i.d. random variables  $X_1, \dots, X_n$  from  $F$  with  $F(\theta) = 1$  for some  $\theta$ , and let  $X_{(m)}^*$  be the maximum of  $X_{(1)}^*, \dots, X_{(m)}^*$  which are i.i.d. from the empirical distribution  $F_n$ .
- Although  $X_{(n)} \rightarrow \theta$ , it never equals  $\theta$ . But  $P_*\{X_{(n)}^* = X_{(n)}\} = 1 - (1 - n^{-1})^n \rightarrow 1 - e^{-1}$ , which leads to the inconsistency of the bootstrap estimator.
- The reason for the inconsistency of the bootstrap is that the bootstrap samples are drawn from  $F_n$  which is not exactly  $F$ . Therefore, the bootstrap may fail due to the lack of “continuity.”

Consider the following problem.

- Let  $X_1, \dots, X_n$  be independent, with a common  $N(\mu, \sigma^2)$  distribution.

- Let  $s^2$  be the sample variance.
- Consider the pivot  $s^2/\sigma^2$ .  
This is distributed, of course, as  $\chi_{n-1}^2/(n-1)$ .
- Consider the bootstrap approximation, namely the distribution of  $s^{*2}/s^2$ , where  $s^{*2}$  is the variance of the resampled data.
- For  $n = 20$ , the bootstrap approximation is not good.
- One source of this problem is in the tails of the normal: about 30% of the variance is contributed by 5% of the distribution and will be missed by typical samples of size 20.
- In other words,  $s^2$  is mean unbiased for  $\sigma^2$ , but quite skewed for moderate  $n$ : in most samples,  $s^2$  is somewhat too small, counterbalanced by the few samples where  $s^2$  is huge. The variance of  $s^{*2}/s^2$  is largely controlled by the sample fourth moment, which is even more skewed.
- For large  $n$ , these problems go away: after all,  $s^2$  and the bootstrap are consistent.

**Project 1.5** Use simulation to illustrate the points made in the above two questions. For the problem of estimating the end point, you can assume that  $X$  is a uniform random variable over  $[0, 1]$  (i.e.,  $\theta = 1$ ).

## Bias Reduction via the Bootstrap Principle

- The bootstrap can also be used to estimate bias and do bias reduction.
- Consider  $\theta_0 = \theta(F_0) = \mu^3$ , where  $\mu = \int x dF_0(x)$ . Set  $\hat{\theta} = \theta(F_n) = \bar{X}^3$ .

- Elementary calculations show that

$$\begin{aligned} E\{\theta(F_n)|F_0\} &= E\left\{\mu + n^{-1} \sum_{i=1}^n (X_i - \mu)\right\}^3 \\ &= \mu^3 + n^{-1} 3\mu\sigma^2 + n^{-2}\gamma, \end{aligned}$$

where  $\gamma = E(X_1 - \mu)^3$  denotes population skewness.

- Using the nonparametric bootstrap, we obtain the following:

$$E\{\theta(F_n^*)|F_n\} = \bar{X}^3 + n^{-1} 3\bar{X}\hat{\sigma}^2 + n^{-2}\hat{\gamma},$$

where  $\hat{\sigma}^2 = n^{-1} \sum (X_i - \bar{X})^2$  and  $\hat{\gamma} = n^{-1} \sum (X_i - \bar{X})^3$  denote sample variance and skewness respectively.

- Using the bootstrap principle,  $E\{\theta(F_n^*)|F_n\} - \theta(F_n)$  is used to estimate  $\theta(F_n) - \theta(F_0)$ .
- Note that  $\theta_0 = \theta(F_n) - (\theta(F_n) - \theta_0)$ . Or,  $\theta_0$  can be estimated by  $\theta(F_n) - [E\{\theta(F_n^*)|F_n\} - \theta(F_n)]$  or  $2\theta(F_n) - E\{\theta(F_n^*)|F_n\}$ .

- The bootstrap bias-reduced estimate is  $2\bar{X}^3 - (\bar{X}^3 + n^{-1}3\bar{X}\hat{\sigma}^2 + n^{-2}\hat{\gamma})$ . Or,  $\hat{\theta}_{NB} = \bar{X}^3 - n^{-1}3\bar{X}\hat{\sigma}^2 - n^{-2}\hat{\gamma}$ .