

STAT575
Large Sample Theory

Pingbang Hu

January 18, 2024

Abstract

This is a graduate-level theoretical statistics course taught by [Georgios Fellouris](#) at University of Illinois Urbana-Champaign, aiming to provide an introduction to asymptotic analysis of various statistical methods, including weak convergence, Lindeberg-Feller CLT, asymptotic relative efficiency, etc.

We list some references of this course, although we will not follow any particular book page by page: *Asymptotic Statistics* [[Vaa98](#)], *Asymptotic Theory of Statistics and Probability* [[Das08](#)], *A course in Large Sample Theory* [[Fer17](#)], *Approximation Theorems of Mathematical Statistics* [[Ser09](#)], and *Elements of Large-Sample Theory* [[Leh04](#)].



This course is taken in Spring 2024, and the date on the cover page is the last updated time.

Contents

| | | |
|----------|--|----------|
| 1 | Introduction | 2 |
| 1.1 | Parametrized Approach | 2 |
| 1.2 | Hypothesis Testing | 2 |
| 1.3 | Different Modes of Convergence | 3 |

Chapter 1

Introduction

Lecture 1: Introduction to Large Sample Theory

Say we first collect n data points $x_1, \dots, x_n \in \mathbb{R}^d$, large sample theory concerns with the limiting theory as $n \rightarrow \infty$. We may treat x_i as a realization of a random vector X_i on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$. In this course, we will primarily consider the case that X_i 's are i.i.d., i.e., independent and identically distributed from a distribution function, or the *cumulative density function* (CDF) F such that

16 Jan. 9:30

$$X = (X^1, \dots, X^d) \sim F(x_1, \dots, x_d) \equiv \mathbb{P}(X^1 \leq x_1, \dots, X^d \leq x_d)$$

for all $x_i \in \mathbb{R}$. If we have access to F , we can compute the corresponding *probability density function* (PDF) \mathbb{P} , and then have access to $\mathbb{P}(X \in A)$ for all (measurable) $A \subseteq \mathbb{R}^d$ of interest. If we know any of the above, we know every thing about the population. Hence, the goal is to compute this by collecting data x_i 's, which is a statistical inference problem.

1.1 Parametrized Approach

There are various ways of doing this task, one way is the so-called parametrized approach. By postulating a family of CDFs $\{F_\theta, \theta \in \Theta\}$ where Θ is often a subset of \mathbb{R}^m for some m (generally $\neq n$), the goal is to select a member of this family that is the “closest”, or the “best fit” to the truth, i.e., F , based on the data. To emphasize that this depends on the data, we sometimes write the function we found as $\hat{\theta}_n(x_1, \dots, x_n)$ so that $F_{\hat{\theta}_n(x_1, \dots, x_n)}$ is our proxy for F .

Now, assume that the family is initially given, the problem is then how to select $\hat{\theta}_n$. Fisher suggested that we should look at the maximum likelihood estimator (MLE). The justification for MLE is not about finite n , but about its asymptotic behavior when $n \rightarrow \infty$. Specifically, we have the following theorem due to Fisher (informally stated).

Theorem 1.1.1 (Fisher). If $F \in \{F_\theta: \theta \in \Theta\}$, i.e., if $F = F_{\theta^*}$ for some $\theta^* \in \Theta$, then under certain conditions, $\hat{\theta}_n$ will be “close” to θ^* as $n \rightarrow \infty$. Under some other conditions, $\sqrt{n}(\hat{\theta}_n - \theta)$ is approximately Gaussian with variance being the “best possible” in some sense.

On the other hand, in the misspecified case, i.e., $F \notin \{F_\theta, \theta \in \Theta\}$, we can still compute the MLE, which leads to another justification for MLE since even in this case, $\hat{\theta}_n$ will still be “close” to θ^* such that F_{θ^*} is, in some sense, the “closest” to F among all possible F_θ (minimizing divergence, to be precise).

1.2 Hypothesis Testing

We will also develop theory for hypothesis testing for some hypothesis we're interested in, e.g., whether the data we collect is really i.i.d., or whether our proposed family is reasonable enough. Say now X_i 's are scalar random variable with $\mathbb{E}[X] = \mu$, and we want to test the null hypothesis $H_0: \mu = 0$.

Example. Consider a controlled group Z and a treatment group Y , and we observe Z_1, \dots, Z_n , and Y_1, \dots, Y_n , respectively, and compute $X_i = Z_i - Y_i$ for all i . Testing H_0 on the distribution of X will show the effect of the treatment.

To do this, a well-known method is the so-called t -test. Let s_n to be the sample standard derivation, then we can compute

$$T_n = \frac{\bar{X}_n}{s_n/\sqrt{n}} \sim t_{n-1}$$

as long as X is Gaussian, i.e., the t -statistics for H_0 . What if X is not an Gaussian? We will show that even if X is not Gaussian, this result is “approximately valid” when n is “large enough” as long as $\text{Var}[X] < \infty$.

Remark (Sample Size). When we say n is “large enough”, what we mean really depends on how fast the underlying distribution will approach Gaussian as n grows. Hence, if we can say more about the underlying population, we can say more about when does n is “large enough”; otherwise such a limiting theory might be completely useless in practice.

What if now $\text{Var}[X]$ doesn’t exit? When the population has a heavy tail distribution, then second moment may not exit.

Example (Cauchy distribution). The Cauchy distribution doesn’t have finite moment of order greater than 1.

In this case, some other test is needed. A simple test would be looking at the sign of X_i , i.e., the sign test.

Example (Sign test). We might reject H_0 if $\sum_{i=1}^n \mathbb{1}_{X_i > 0}$ is large. Note that under H_0 , $\sum_{i=1}^n \mathbb{1}_{X_i > 0} \sim \text{Bin}(n, 1/2)$, and this test is valid even if expectation doesn’t exist.

We see that without saying anything about F , the sign test is valid even for $n = 3$ or 5 as the sum is exactly binomial distribution under H_0 . Although simple and have good property, only looking at the sign of X_i might be too weak. A natural idea is to look at the absolute value of X_i .

Example (Wilcoxon’s rank-sum test). Let $R_{i,n}$ to be the rank of $|X_i|$, then consider the so-called *Wilcoxon’s rank-sum test*

$$\sum_{i=1}^n \mathbb{1}_{X_i > 0} R_{i,n}.$$

As one can imagine, the closed form of the above sum will be complicated; however, asymptotically, the above statics will follow Gaussian again, such that the rate of convergence doesn’t depend on the underlying population.

Finally, we also ask how can we compare these different tests? This will also be addressed in this course.

Lecture 2: Modes of Convergence

1.3 Different Modes of Convergence

18 Jan. 9:30

Given a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, consider a sequence of d -dimensional random vectors (X_n) and a random vector X , i.e., $X_n, X: \Omega \rightarrow \mathbb{R}^d$. We now discuss different modes of convergence for (X_n) .

Definition 1.3.1 (Point-wise convergence). (X_n) *point-wise converges* to X , denoted as $X_n \rightarrow X$, if $X_n(\omega) \rightarrow X(\omega)$ for all $\omega \in \Omega$.

As previously seen. From analysis, $X_n(\omega) \rightarrow X(\omega)$ if and only if for every $\epsilon > 0$, there exists

$n_0(\omega) \in \mathbb{N}$ such that for every $n \geq n_0$, $\|X_n(\omega) - X(\omega)\|_2 < \epsilon$.

However, since we don't care about measure zero sets, we may instead consider the following.

Definition 1.3.2 (Almost-surely convergence). (X_n) *almost-surely converges* to X , denoted as $X_n \xrightarrow{\text{a.s.}} X$, if $\mathbb{P}(X_n \rightarrow X) = 1$.

In other words, **almost-surely convergence** means that $X_n(\omega) \rightarrow X(\omega)$ for all $\omega \in \Omega \setminus N$ where $\mathbb{P}(N) = 0$. However, this might still be too strong.

Definition 1.3.3 (Convergence in probability). (X_n) *converges in probability* to X , denoted as $X_n \xrightarrow{p} X$, if for every $\epsilon > 0$, $\mathbb{P}(\|X_n - X\| > \epsilon) \rightarrow 0$ as $n \rightarrow \infty$.

Remark. $X_n \rightarrow X$ if and only if $\|X_n - X\| \rightarrow 0$. The same also holds for \xrightarrow{p} and $\xrightarrow{\text{a.s.}}$.

A related notion is the following, where we now sum over n .

Definition 1.3.4 (Converges completely). (X_n) *converges completely* to X , denoted as $X_n \xrightarrow{\text{comp}} X$, if for every $\epsilon > 0$, $\sum_{n=1}^{\infty} \mathbb{P}(\|X_n - X\| > \epsilon) < \infty$.

Finally, we have the following.

Definition 1.3.5 (Converges in L^p). Let $p > 0$, we say $X_n \xrightarrow{L^p} X$ if $\mathbb{E}[\|X_n - X\|^p] \rightarrow 0$ as $n \rightarrow \infty$.

1.3.1 Connection Between Modes of Convergence

We have the following connections between different modes of convergence.

$$\text{completely} \implies \text{almost-surely} \implies \text{in probability} \longleftarrow \text{in } L^p$$

To show the above, the following characterization for **almost-surely convergence** is useful.

Proposition 1.3.1. For a sequence of random vectors (X_n) and a random vector X , we have

$$\begin{aligned} X_n \xrightarrow{\text{a.s.}} X &\Leftrightarrow \mathbb{P}(\|X_k - X\| > \epsilon \text{ for some } k \geq n) \xrightarrow{n \rightarrow \infty} 0 \\ &\Leftrightarrow \mathbb{P}(\|X_n - X\| > \epsilon \text{ for infinitely many } n\text{'s}) = 0 \\ &\Leftrightarrow \mathbb{P}(\limsup_{n \rightarrow \infty} \|X_n - X\| > \epsilon) = 0, \end{aligned}$$

where the above holds for every $\epsilon > 0$.

From **Proposition 1.3.1**, it's clear that $\xrightarrow{\text{a.s.}}$ implies \xrightarrow{p} since

$$\mathbb{P}(\|X_k - X\| > \epsilon \text{ for some } k \geq n) \geq \mathbb{P}(\|X_n - X\| > \epsilon),$$

hence if the former goes to 0, so does the latter. On the other hand, $\xrightarrow{\text{comp}}$ implies $\xrightarrow{\text{a.s.}}$ follows from the third equivalence. Lastly, the **convergence in L^p** implies the **convergence in probability** since

$$\mathbb{P}(\|X_n - X\| > \epsilon) \leq \frac{1}{\epsilon^p} \mathbb{E}[\|X_n - X\|^p]$$

from Markov's inequality. However, the converse is not always true.

Theorem 1.3.1 (Dominated convergence theorem). If $X_n \xrightarrow{p} X$ and $\|X_n - X\| \leq Z$ for all $n \geq 1$ where $\mathbb{E}[\|Z\|^p] < \infty$, then $X_n \xrightarrow{L^p} X$.

Theorem 1.3.2 (Scheffé's theorem). If $X_n \xrightarrow{p} X$ and $\limsup_{n \rightarrow \infty} \mathbb{E} [\|X_n\|^p] \leq \mathbb{E} [\|X\|^p] < \infty$, then $X_n \xrightarrow{L^p} X$.

1.3.2 Applications to Statistics

Let $(X_n) \stackrel{\text{i.i.d.}}{\sim} F$ where F is a distribution function. Say we're interested in some aspect of F , for example, some parameter $\theta = T(F) \in \mathbb{R}^m$. By collecting data X_1, \dots, X_n , we estimate θ by computing an estimator $\hat{\theta}_n$ of θ .¹

Definition 1.3.6 (Consistent). $\hat{\theta}_n$ is *consistent* of θ if $\hat{\theta}_n \xrightarrow{p} \theta$ as $n \rightarrow \infty$.

Definition 1.3.7 (Strongly consistent). $\hat{\theta}_n$ is *strongly consistent* of θ if $\hat{\theta}_n \xrightarrow{\text{a.s.}} \theta$ as $n \rightarrow \infty$.

Let's first see the most well-known estimation problem, the mean estimation.

Example (Mean estimation). Suppose $d = 1$, and let X be non-negative. Say we're interested in $\theta = \mathbb{E}[X]$. It's standard that in this case, we can compute $\mathbb{E}[X]$ by

$$\theta = \mathbb{E}[X] = \int_0^\infty \mathbb{P}(X > t) dt = \int_0^\infty (1 - F(t)) dt.$$

On the other hand, if X has a PMF f , then

$$\mathbb{E}[X] = \sum_x x f(x) = \sum_x x \Delta F(x),$$

where $f(x) = \Delta F(x) \equiv F(x) - F(x_-)$. And if X has a PDF f , then

$$\mathbb{E}[X] = \int_0^\infty x f(x) dx = \int_0^\infty x F(dx)$$

where $F(dx) := f(x)dx$ in a measure-theoretical sense.

Now, let $\hat{\theta}_n$ to be the sample mean, i.e., $\hat{\theta}_n = \bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$. From the strong law of large number, $\bar{X}_n \xrightarrow{\text{a.s.}} \mathbb{E}[X]$, which implies that $\hat{\theta}_n$ is a **strongly consistent estimator** of θ .

On the other hand, if $\text{Var}[X] < \infty$, then $\bar{X}_n \xrightarrow{L^2} \mathbb{E}[X]$, which further implies $\bar{X}_n \xrightarrow{p} \mathbb{E}[X]$, hence $\hat{\theta}_n$ is **consistent**.^a

^aThe latter is true even without $\text{Var}[X] < \infty$ as we expect.

Proof. We show the last statement. Since $\text{Var}[X] < \infty$, then

$$\frac{\text{Var}[X]}{n} = \text{Var}[\bar{X}_n] = \mathbb{E}[(\bar{X}_n - \mathbb{E}[X])^2] \rightarrow 0$$

as $n \rightarrow \infty$, which implies $\bar{X}_n \xrightarrow{p} \mathbb{E}[X]$. ⊛

Another interesting problem is the supremum estimation.

Example (Supremum estimation). Suppose there is a $\theta \in \mathbb{R}$ where distribution function F such that $F(\theta - \epsilon) < 1 = F(\theta)$ for all $\epsilon > 0$. This means $\theta = \sup_\omega X(\omega)$ since $\mathbb{P}(X \leq \theta - \epsilon) = F(\theta - \epsilon)$ and $F(\theta) = \mathbb{P}(X \leq \theta)$.^a The natural estimator for θ would be $\hat{\theta}_n = \max_{1 \leq i \leq n} X_i$, and it's indeed **strongly consistent**.

^aSuch an distribution exists, for example, $\mathcal{U}(0, \theta)$.

¹ $\hat{\theta}_n$ is a function of X_i 's.

Proof. We see that for any $\epsilon > 0$,

$$\begin{aligned}\mathbb{P}(|\hat{\theta}_n - \theta| > \epsilon) &= \mathbb{P}(\hat{\theta}_n > \theta + \epsilon) + \mathbb{P}(\hat{\theta}_n < \theta - \epsilon) \\ &= \mathbb{P}\left(\bigcup_{i=1}^n \{X_i > \theta + \epsilon\}\right) + \mathbb{P}\left(\bigcap_{i=1}^n \{X_i < \theta - \epsilon\}\right) \\ &\leq \sum_{i=1}^n \underbrace{\mathbb{P}(X > \theta + \epsilon)}_0 + \prod_{i=1}^n \mathbb{P}(X_i < \theta - \epsilon) = (\mathbb{P}(X_1 < \theta - \epsilon))^n \leq (F(\theta - \epsilon))^n \rightarrow 0\end{aligned}$$

as $n \rightarrow \infty$ since $F(\theta - \epsilon) < 1$. This shows that $\hat{\theta}_n$ is indeed **consistent**. Moreover, since $\mathbb{P}(|\hat{\theta}_n - \theta| > \epsilon)$ decays exponentially, so this is absolutely summable, hence it's also **strongly consistency**. \circledast

Proving convergence of $\hat{\theta}_n$ is useful, but this might not be enough.

Example. Consider any deterministic sequence (a_n) in \mathbb{R} which converges to 0. Adding a_n to $\hat{\theta}_n$ will not change the convergence of $\hat{\theta}_n$. This shows that being **consistent** might not be enough in some cases.

The above suggests that we should look at the *distribution* of $\hat{\theta}_n - \theta$ in order to say how does $\hat{\theta}_n \rightarrow \theta$.

Example (Mean estimation for Gaussian). Suppose $X \sim \mathcal{N}(\theta, 1)$. Then $\hat{\theta}_n = \bar{X}_n \sim \mathcal{N}(\theta, 1/n)$, i.e., $\sqrt{n}(\hat{\theta}_n - \theta) \sim \mathcal{N}(0, 1)$. This implies that we can write down a confidence interval (CI) such that $\hat{\theta}_n \pm 1.96/\sqrt{n}$ with 95% CI for $\hat{\theta}_n$.

Doing this for other kind of estimators and F is not that straightforward and will be challenging.

Remark. Let (X_n) and X be d -dimensional random vectors, $h: \mathbb{R}^d \rightarrow \mathbb{R}^m$, and $c \in \mathbb{R}^d$ constant.

- (a) If $X_n \rightarrow c$, then $h(X_n) \rightarrow h(c)$ if h is continuous at c .^a This also holds for $\xrightarrow{\text{a.s.}}$ and \xrightarrow{P} .
- (b) If $X_n \rightarrow X$, then $h(X_n) \rightarrow h(X)$ if h is continuous. This also holds for $\xrightarrow{\text{a.s.}}$ and \xrightarrow{P} .

^aThis is an if and only if condition if this holds for any h .

Let's see some examples.

Example. If $d = 1$, and $X_n \rightarrow \theta \neq 0$. Then $1/X_n \rightarrow 1/\theta$ where

$$h(x) = \begin{cases} \frac{1}{x}, & \text{if } x \neq 0; \\ c, & \text{if } x = 0 \end{cases}$$

for any $c \in \mathbb{R}$. The same holds for $\xrightarrow{\text{a.s.}}$ and \xrightarrow{P} .

Example. If $X_n \rightarrow X$ and $Y_n \rightarrow Y$, then $(X_n Y_n) \rightarrow (X, Y)$.^a The same holds for $\xrightarrow{\text{a.s.}}$ and \xrightarrow{P} .

^aThe converse is also true since projections are continuous.

Proof. j To show $\|(X_n, Y_n) - (X, Y)\| \rightarrow 0$, we have

$$\|(X_n, Y_n) - (X, Y)\| \leq \|X_n - X\| + \|Y_n - Y\|$$

for all $n \geq 1$.^a The latter two terms goes to 0 (in whatever sense) by assumption. \circledast

^aThis can be seen from $\sqrt{x+y} \leq \sqrt{x} + \sqrt{y}$.

Appendix

Bibliography

- [Das08] Anirban DasGupta. *Asymptotic Theory of Statistics and Probability*. Springer Science & Business Media, Feb. 6, 2008. 727 pp. ISBN: 978-0-387-75971-5. Google Books: [sX4_AAAAQBAJ](#).
- [Fer17] Thomas S. Ferguson. *A Course in Large Sample Theory*. Routledge, Sept. 6, 2017. 140 pp. ISBN: 978-1-351-47005-6. Google Books: [clcODwAAQBAJ](#).
- [Leh04] E. L. Lehmann. *Elements of Large-Sample Theory*. Springer Science & Business Media, Aug. 27, 2004. 640 pp. ISBN: 978-0-387-98595-4. Google Books: [geIoxvgTXlEC](#).
- [Ser09] Robert J. Serfling. *Approximation Theorems of Mathematical Statistics*. John Wiley & Sons, Sept. 25, 2009. 399 pp. ISBN: 978-0-470-31719-8. Google Books: [enUouJ4EHzQC](#).
- [Vaa98] A. W. van der Vaart. *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge: Cambridge University Press, 1998. ISBN: 978-0-521-78450-4. DOI: [10.1017/CB09780511802256](#). URL: <https://www.cambridge.org/core/books/asymptotic-statistics/A3C7DAD3F7E66A1FA60E9C8FE132EE1D> (visited on 10/17/2023).