STAT575 Lrage Sample Theory

Pingbang Hu

February 3, 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

	Introduction
	1.1 Parametrized Approach
	1.2 Hypothesis Testing
2	Modes of Convergence
	2.1 Different Modes of Convergence
	2.2 Weak Convergence
	2.3 Stochatsic Boundedness

Chapter 1

Introduction

Lecture 1: Introduction to Large Sample Theory

Say we first collect n data points $x_1, \ldots, x_n \in \mathbb{R}^d$, large sample theory concerns with the limiting theory as $n \to \infty$. We may treat x_i as a realization of a random vector X_i on a probability space $(\Omega, \mathscr{F}, \mathbb{P})$. 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

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

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.

Remark. 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 "closet", or the "best fit" to the truth, i.e., F, based on the data.

Note. To emphasize that this depends on the data, we sometimes write the function we found as $\hat{\theta}_n(x_1,\ldots,x_n)$ so that $F_{\hat{\theta}_n(x_1,\ldots,x_n)}$ is our proxy for F.

Now, assume that the family is initially given, the problem is then how to select $\hat{\theta}_n$.

Example. 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 $n \to \infty$. 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 $n \to \infty$. 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, \ldots, Z_n , and Y_1, \ldots, 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

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

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 $\operatorname{Var}[X] < \infty$.

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 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.

Chapter 2

Modes of Convergence

Lecture 2: Modes of Convergence

2.1 Different Modes of Convergence

18 Jan. 9:30

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

Definition 2.1.1 (Point-wise converge). (X_n) point-wise converges to X, denoted as $X_n \to X$, if $X_n(\omega) \to X(\omega)$ for all $\omega \in \Omega$.

 ${}^a\mathrm{I.e.}, \text{ for every } \epsilon > 0, \text{ there exists } n_0(\omega) \in \mathbb{N} \text{ such that for every } n \geq n_0, \|X_n(\omega) - X(\omega)\|_2 < \epsilon.$

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

Definition 2.1.2 (Almost-surely converge). (X_n) almost-surely converges to X, denoted as $X_n \stackrel{\text{a.s.}}{\to} X$, if $\mathbb{P}(X_n \to X) = 1$.

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

Definition 2.1.3 (Converge in probability). (X_n) converges in probability to X, denoted as $X_n \stackrel{p}{\to} X$, if for every $\epsilon > 0$, $\mathbb{P}(||X_n - X|| > \epsilon) \to 0$ as $n \to \infty$.

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

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

Definition 2.1.4 (Converge completely). (X_n) converges completely to X, denoted as $X_n \stackrel{\text{comp}}{\to} X$, if for every $\epsilon > 0$, $\sum_{n=1}^{\infty} \mathbb{P}(\|X_n - X\| > \epsilon) < \infty$.

Finally, we have the following.

Definition 2.1.5 (Converge in L^p). (X_n) converges in L^p to X for some p > 0, denoted as $X_n \stackrel{L^p}{\to} X$, if $\mathbb{E}[\|X_n - X\|^p] \to 0$ as $n \to \infty$.

2.1.1 Connection Between Modes of Convergence

We have the following connections between different modes of convergence.

completely \Longrightarrow almost-surely \Longrightarrow in probability \Longleftrightarrow in L^p

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

Proposition 2.1.1. For a sequence of random vectors (X_n) and a random vector X, we have

$$X_n \stackrel{\text{a.s.}}{\to} X \Leftrightarrow \mathbb{P}(\|X_k - X\| > \epsilon \text{ for some } k \ge n) \stackrel{n \to \infty}{\to} 0$$

 $\Leftrightarrow \mathbb{P}(\|X_n - X\| > \epsilon \text{ for infinitely many } n's) = 0$
 $\Leftrightarrow \mathbb{P}(\limsup_{n \to \infty} \|X_n - X\| > \epsilon) = 0,$

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

From Proposition 2.1.1, it's clear that $\stackrel{\text{a.s.}}{\rightarrow}$ implies $\stackrel{p}{\rightarrow}$ since

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

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

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

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

Theorem 2.1.1 (Dominated convergence theorem). If $X_n \stackrel{p}{\to} X$ and $||X_n - X|| \le Z$ for all $n \ge 1$ where $\mathbb{E}[||Z||^p] < \infty$, then $X_n \stackrel{L^p}{\to} X$.

Theorem 2.1.2 (Scheffé's theorem). If $X_n \stackrel{p}{\to} X$ and $\limsup_{n \to \infty} \mathbb{E}\left[\|X_n\|^p\right] \le \mathbb{E}\left[\|X\|^p\right] < \infty$, then $X_n \stackrel{L^p}{\to} X$.

2.1.2 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, \ldots, X_n , we estimate θ by computing an estimator $\hat{\theta}_n$ of θ .¹

Definition 2.1.6 (Consistent). $\hat{\theta}_n$ is consistent of θ if $\hat{\theta}_n \stackrel{p}{\to} \theta$ as $n \to \infty$.

Definition 2.1.7 (Strongly consistent). $\hat{\theta}_n$ is strongly consistent of θ if $\hat{\theta}_n \stackrel{\text{a.s.}}{\to} \theta$ as $n \to \infty$.

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

Example (Mean esimation). 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

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

 $^{{}^{1}\}hat{\theta}_{n}$ is a function of X_{i} 's.

*

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

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

On the other hand, if $\operatorname{Var}[X] < \infty$, then $\overline{X}_n \stackrel{L^2}{\to} \mathbb{E}[X]$, which further implies $\overline{X}_n \stackrel{p}{\to} \mathbb{E}[X]$, hence $\hat{\theta}_n$ is consistent.

^aThe latter is true even without $\operatorname{Var}\left[X\right]=\infty$ as we expect.

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

$$\frac{\operatorname{Var}\left[X\right]}{n} = \operatorname{Var}\left[\overline{X}_{n}\right] = \mathbb{E}\left[\left(\overline{X} - \mathbb{E}\left[X\right]\right)^{2}\right] \to 0$$

as $n \to \infty$, which implies $\overline{X}_n \stackrel{p}{\to} \mathbb{E}[X]$.

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 $F(\theta) = \mathbb{P}(X \leq \theta)$. 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, $\mathcal{U}(0,\theta)$.

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

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

as $n \to \infty$ 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.

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 $\hat{\theta}_n \to \theta$.

Example (Mean estimation for Gaussian). Suppose $X \sim \mathcal{N}(\theta, 1)$. Then $\hat{\theta}_n = \overline{X}_n \sim \mathcal{N}(\theta, 1/n)$, 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, $h: \mathbb{R}^d \to \mathbb{R}^m$, and $c \in \mathbb{R}^d$ constant.

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

Let's see some examples.

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

Example. If d=1, and $X_n \to \theta \neq 0$. Then $1/X_n \to 1/\theta$ 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 $\overset{\text{a.s.}}{\to}$ and $\overset{p}{\to}$.

Example. If $X_n \to X$ and $Y_n \to Y$, then $(X_n Y_n) \to (X,Y)$. The same holds for $\stackrel{\text{a.s.}}{\to}$ and $\stackrel{p}{\to}$.

^aThe converse is also true since projections are continuous.

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

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

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

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

Lecture 3: Weak Convergence Portmanteau theorem

2.2 Weak Convergence

25 Jan. 9:30

All convergences we have discussed are in some senses "point-wise" but not "distribution-wise", and the latter is more powerful. Consider working with a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ and the following.

Definition 2.2.1 (Total variation). The total variation distance between X and Y in Ω is defined as

$$\mathrm{TV}(X,Y) = \sup_{B \in \mathscr{F}} |\mathbb{P}(X \in B) - \mathbb{P}(Y \in B)|$$

Returning to our situation, consider a sequence or random variables (X_n) and a random variable X.

Remark. If X_n has density f_n and X has density f, then $TV(X_n, X) = \frac{1}{2} \int |f_n - f|$.

Definition 2.2.2 (Convergence in total variation). (X_n) converges in total variation to X, denoted as $X_n \stackrel{\text{TV}}{\to} X$, if $\text{TV}(X_n, X) \to 0$ as $n \to \infty$.

Remark. If X_n has density f_n and X has density $f, f_n \to f$ implies $X_n \stackrel{\text{TV}}{\to} X$.

Note. The above could make sense even if X_n is defined on different $(\Omega_n, \mathscr{F}_n, \mathbb{P}_n)$ for every n. Let's see some examples.

Let b bee boine examples.

Example. Consider $X_n \sim \text{Bin}(n, p_n)$ such that $np_n \to \lambda \in \mathbb{R}$. As this happens,

$$X_n \sim \text{Bin}(n, p_n) \stackrel{\text{TV}}{\to} X \sim \text{Pois}(\lambda).$$

Example. Let $X_n \sim f_{\theta_n}$ where $f_{\theta_n}(x) = f(x)e^{\theta x - \psi(\theta)}$ for some $\theta \in \Theta$. If $\theta_n \to \theta$, then $X_n \stackrel{\mathrm{TV}}{\to} X \sim f_{\theta}$. For example, if $X_n \sim \mathrm{Pois}(\theta_n)$ and $\theta_n \to \theta$, then $X_n \stackrel{\mathrm{TV}}{\to} X \sim \mathrm{Pois}(\theta)$.

However, convergence in total variation might be too strong to work with.

Example. Let $X_n \sim \mathcal{U}\{0, 1/n, \dots, (n-1)/n\}$, which should be converging to $X \sim \mathcal{U}(0, 1)$. However, this doesn't happen in total variation distance as we can take B to be \mathbb{Q} .

This suggests that we should look at something weaker.

Definition 2.2.3 (Weak convergence). (X_n) converges weakly to X, denoted as $X_n \stackrel{\text{w}}{\to} X$, if for all bounded continuous $g: \mathbb{R}^d \to \mathbb{R}$, $\mathbb{E}[g(X_n)] \to \mathbb{E}[g(X)]$.

To see how is weak convergence compared to convergence in total variation, we revisit the above.

Example. Let $X_n \sim \mathcal{U}\{0, 1/n, \dots, (n-1)/n\}$, which should be converging to $X \sim \mathcal{U}(0, 1)$. We have

$$\mathbb{E}\left[g(X_n)\right] = \sum_{k=0}^{n-1} g(k/n) \left(\frac{k+1}{n} - \frac{k}{n}\right) \to \int_0^1 g(x) \, \mathrm{d}x = \mathbb{E}\left[g(X)\right]$$

as g is bounded and continuous on [0,1], hence Riemann integrable.

2.2.1 Portmanteau Theorem

The following is our main tool of proving weak convergence.

Theorem 2.2.1 (Portmanteau theorem). The following are equivalent.

- (a) $X_n \stackrel{\text{w}}{\to} X$.
- (b) $\mathbb{E}\left[g(X_n)\right] \to \mathbb{E}\left[g(X)\right]$ for all bounded Lipschitz $g \colon \mathbb{R}^d \to \mathbb{R}$.
- (c) $\mathbb{P}(X \in A) \leq \liminf_{n \to \infty} \mathbb{P}(X_n \in A)$ for all $A \subseteq \mathbb{R}^d$ open.
- (d) $\mathbb{P}(X \in A) \ge \limsup_{n \to \infty} \mathbb{P}(X_n \in A)$ for all $A \subseteq \mathbb{R}^d$ closed.
- (e) $\mathbb{P}(X_n \in A) \to \mathbb{P}(X \in A)$ for all A such that $\mathbb{P}(X \in \partial A) = 0$.

Before we prove Portmanteau theorem, we should note that all our discussion can be extended to metric spaces from Euclidean spaces. Let's first recall some basic results for metric spaces.

Claim. Given a metric space (S, ρ) , $\rho(\cdot, A)$ is Lipschitz for any $A \subseteq S$, i.e., for any $x, y \in S$,

$$|\rho(x, A) - \rho(y, A)| \le \rho(x, y).$$

Proof. Since for any $z \in S$, $\rho(x,z) \le \rho(x,y) + \rho(y,z)$, hence $\rho(x,A) - \rho(y,A) \le \rho(x,y)$ by taking the infimum over $z \in A$. Interchanging x and y gives another inequality.

Claim. Given a metric space (S, ρ) , for any $A \subseteq S$, $x \in \overline{A} \Leftrightarrow \rho(x, A) = 0$.

Proof. If $x \in \overline{A}$, there exists (x_n) in A such that $\rho(x_n, x) \to 0$. Then for any $z \in A$, $\rho(x, z) \le \rho(x, x_n) + \rho(x_n, z)$, implying

$$\rho(x,A) < \rho(x,x_n) + \rho(x_n,A) \to 0$$

hence $\rho(x,A)=0$. On the other hand, suppose $\rho(x,A)=0$. As $\rho(x,A)=\inf_{y\in A}\rho(x,y)$, there exists (y_n) in A such that $\rho(x,y_n)\to\rho(x,A)=0$, i.e., $x\in\overline{A}$.

The crucial lemma we're going to use to prove Portmanteau theorem is the following.

Lemma 2.2.1. Given a metric space (S, ρ) and let $A \subseteq S$ be a closed subset. Then there exists bounded Lipschitz $g_k \colon S \to \mathbb{R}$, decreasing in k such that $g_k(x) \searrow \mathbb{1}_A(x)$.

Proof. Since A is closed, $A = \overline{A}$ and

$$\mathbb{1}_{A}(x) = \begin{cases} 1, & \text{if } x \in A \Leftrightarrow \rho(x, A) = 0; \\ 0, & \text{if } x \notin A \Leftrightarrow \rho(x, A) > 0. \end{cases}$$

Now, we let

$$g_k(x) = \begin{cases} 0, & \text{if } \rho(x, A) > \frac{1}{k}; \\ 1 - k\rho(x, A), & \text{otherwise;} \end{cases} = 1 - (k\rho(x, A) \wedge 1).$$

We see that

- if $x \in A$: $\mathbb{1}_A(x) = 1$, and $g_k(x) = 1$ since $\rho(x, A) = 0$;
- if $x \notin A$: $\mathbb{1}_A(x) = 0$, and $\rho(x, A) > 0$ since A closed, and $g_k(x) = 0$ for all large enough k.

Finally, it's clear that $g_k(x)$ takes values in [0, 1], and we now show it's Lipschitz. We have

$$|g_k(x) - g_k(y)| = |(k\rho(x, A) \wedge 1) - (k\rho(y, A) \wedge 1)| \le k\rho(x, y)$$

for all $x, y \in S$.

Then we can prove the Portmanteau theorem.

Proof of Theorem 2.2.1. (a) \Rightarrow (b) is clear. And we start by proving (c) \Leftrightarrow (d).

Claim. (c) \Leftrightarrow (d).

Proof. We first prove that $(d) \Rightarrow (c)$. Since when A is open,

$$\mathbb{P}(X \in A) = 1 - \mathbb{P}(X \in A^c)
\leq 1 - \limsup_{n \to \infty} \mathbb{P}(X_n \in A^c)
= 1 - \limsup_{n \to \infty} (1 - \mathbb{P}(X_n \in A)) = \liminf_{n \to \infty} \mathbb{P}(X_n \in A).$$
(d)

$$(c) \Rightarrow (d)$$
 is exactly the same, hence $(c) \Leftrightarrow (d)$.

Next, we prove (b) \Rightarrow (d), which gives us (a) \Rightarrow (b) \Rightarrow (d) \Leftrightarrow (c).

Claim. (b) \Rightarrow (d).

Proof. From Lemma 2.2.1, there exists bounded Lipschitz $g_k \searrow \mathbb{1}_A$ such that for all closed A,

$$\mathbb{P}(X_n \in A) = \mathbb{E}\left[\mathbb{1}_A(X_n)\right] \leq \mathbb{E}\left[g_k(X_n)\right].$$

This is true for every k and n since $g_k \geq \mathbb{1}_A$, and by taking the limit as $n \to \infty$,

$$\lim\sup_{n\to\infty} \mathbb{P}(X_n\in A) \leq \lim\sup_{n\to\infty} \mathbb{E}\left[g_k(X_n)\right] = \mathbb{E}\left[g_k(X)\right]$$

from our assumption (b). Finally, as $k \to \infty$, it goes to $\mathbb{E}[\mathbb{1}_A(X)] = \mathbb{P}(X \in A)$ as desired. \circledast

The proof will be continued...

Lecture 4: Continuous Mapping Theorem

Before finishing the proof of Portmanteau theorem, we need one additional tool.

30 Jan. 9:30

Lemma 2.2.2. If $\{A_i\}_{i\in I}$ are pairwise disjoint events, then $\{i\in I: \mathbb{P}(A_i)>0\}$ is countable.

aNote that I can be uncountable.

Proof. Since we can write

$$\{i \in I \colon \mathbb{P}(A_i) > 0\} = \bigcup_{k=1}^{\infty} \left\{ i \in I \colon \mathbb{P}(A_i) \ge \frac{1}{k} \right\} =: \bigcup_{k=1}^{\infty} I_k,$$

hence it suffices to show $|I_k| < \infty$ for any $k \ge 1$. Indeed, for any k, $|I_k| \le k$. Suppose not. Then there exists a countable $J_k \subseteq I_k$ such that $|J_k| > k$, implying

$$\mathbb{P}\left(\bigcup_{i\in J_k} A_i\right) = \sum_{i\in J_k} \mathbb{P}(A_i) \ge \frac{|J_k|}{k} > 1,$$

which is a contradiction.

We now finish the proof of Portmanteau theorem.

Proof of Theorem 2.2.1 (cont.) We already proved (a) \Rightarrow (b) \Rightarrow (d) \Leftrightarrow (c).

Claim. (c) + (d) \Rightarrow (e).

Proof. We see that for any $A, A^o \subseteq A \subseteq \overline{A}$, and from (c),

$$\mathbb{P}(X \in A^{o}) \leq \liminf_{n \to \infty} \mathbb{P}(X_{n} \in A^{o}) \leq \liminf_{n \to \infty} \mathbb{P}(X_{n} \in A)$$
$$\leq \limsup_{n \to \infty} \mathbb{P}(X_{n} \in A) \leq \limsup_{n \to \infty} \mathbb{P}(X_{n} \in \overline{A}) \leq \mathbb{P}(X \in \overline{A})$$

where the last step follows from (d). Finally, since

$$\mathbb{P}(X \in \overline{A}) - \mathbb{P}(X \in A^o) = \mathbb{P}(\{X \in \overline{A}\} \setminus \{X \in A^o\}) = \mathbb{P}(X \in (\overline{A} \setminus A^o)) = \mathbb{P}(X \in \partial A),$$

which is 0 by our assumption, i.e., inequalities above are all equalities. In particular, since

$$\lim_{n \to \infty} \inf \mathbb{P}(X_n \in A) \le \lim_{n \to \infty} \mathbb{P}(X_n \in A) \le \lim_{n \to \infty} \mathbb{P}(X_n \in A)$$

and
$$\mathbb{P}(X \in A^o) \leq \mathbb{P}(X \in A) \leq \mathbb{P}(X \in \overline{A}), \ \mathbb{P}(X \in A) = \lim_{n \to \infty} \mathbb{P}(X_n \in A).$$

Finally, we prove the following.

Claim. (e) \Rightarrow (a).

Proof. For every $g: \mathbb{R}^d \to \mathbb{R}$ bounded and continuous, we want to show $\mathbb{E}[g(X_n)] \to \mathbb{E}[g(X)]$. Suppose $g \geq 0$, and let $K \geq g(x)$ for every $x \in \mathbb{R}^d$ (which exists since g is bounded), then

$$\mathbb{E}\left[g(X_n)\right] = \int_0^K \mathbb{P}(g(X_n) > t) \, \mathrm{d}t, \quad \mathbb{E}\left[g(X)\right] = \int_0^K \mathbb{P}(g(X) > t) \, \mathrm{d}t,$$

so we just need to prove the convergence of the above two integrals. From bounded convergence theorem, it suffices to show that for almost every $t \in [0, K]$,

$$\mathbb{P}(q(X_n) > t) \to \mathbb{P}(q(X) > t).$$

Observe that $\mathbb{P}(g(X_n) > t) = \mathbb{P}(X_n \in \{g > t\})$ and $\mathbb{P}(g(X) > t) = \mathbb{P}(X \in \{g > t\})$, so from (e) with $A := \{g > t\}$, it suffices to show $\mathbb{P}(X \in \partial \{g > t\}) = 0$ for almost all t. Firstly,

$$\mathbb{P}(X \in \partial \{g > t\}) = \mathbb{P}(X \in \overline{\{g > t\}} \setminus \{g > t\}^o) = \mathbb{P}(X \in \overline{\{g \geq t\}} \setminus \{g > t\}) = \mathbb{P}(g(X) = t).$$

Moreover, consider the events $\{g(X)=t\}_{t\in[0,K]}$, which are pairwise disjoint, hence Lemma 2.2.2 implies $\mathbb{P}(g(X)=t)=0$ for all but countably many t's, exactly what we want to show.

This finishes the proof.

^aOtherwise, we consider $g = g^+ - g^-$ where $g^+ = \max(g, 0)$ and $g^- = \max(-g, 0)$, and everything follows.

2.2.2 Continuous Mapping Theorem

A common scenario is that given a nice function h (in terms of continuity), if $X_n \stackrel{\text{w}}{\to} X$, we want to know when will $h(X_n) \stackrel{\text{w}}{\to} h(X)$. To develop the theorem of this, we need some more facts about metric spaces.

As previously seen. Given two metric spaces (S, ρ) , (S', ρ') , $g: S \to S'$ is continuous if $x_n \stackrel{\rho}{\to} x$ implies $g(x_n) \stackrel{\rho'}{\to} g(x)$, or for open $A \subseteq S'$, $g^{-1}(A) \subseteq S$ is open.

Notation. We sometimes write $g^{-1}(A) =: \{g \in A\}$.

It's clear that the following holds.

Note. If $g: S \to S'$ is continuous and $A \subseteq S'$ is closed, then $\overline{\{g \in A\}} = \{g \in \overline{A}\}.$

However, when g is not continuous and A is not closed, the situation is a bit more complicated. But at least we can first look at the set where g is continuous.

Notation (Continuous set). For any $g: S \to S'$, we denote the *continuous set* as $C_g := \{x \in S : g \text{ is continuous at } x\}$.

Then we have the following.

Proposition 2.2.1. Given $g: S \to S'$ between metric spaces and $A \subseteq S'$,

$$C_g \cap \overline{\{g \in A\}} \subseteq \{g \in \overline{A}\}.$$

Proof. Let $x \in C_g \cap \overline{\{g \in A\}}$. Since $x \in \overline{\{g \in A\}}$, there exists $(x_n) \in \{g \in A\}$ such that $x_n \stackrel{\rho}{\to} x$. Moreover, $x \in C_g$ implies g is continuous at x, hence $g(x_n) \stackrel{\rho'}{\to} g(x)$, i.e., $g(x) \in \overline{A}$.

This allows us to prove the following theorem, which answers our main question in this section.

Theorem 2.2.2 (Continuous mapping theorem). Consider $X_n \stackrel{\text{w}}{\to} X$ and $h: \mathbb{R}^d \to \mathbb{R}^m$. If $\mathbb{P}(X \in C_h) = 1$, then $h(X_n) \stackrel{\text{w}}{\to} h(X)$.

Proof. Let $A \subseteq \mathbb{R}^m$ be a closed set. Then from Portmanteau theorem (d), we need to show

$$\limsup_{n \to \infty} \mathbb{P}(h(X_n) \in A) \le \mathbb{P}(h(X) \in A).$$

Since $\limsup_{n\to\infty} \mathbb{P}(h(X_n)\in A) = \limsup_{n\to\infty} \mathbb{P}(X_n\in\{h\in A\})$, implying

$$\limsup_{n \to \infty} \mathbb{P}(h(X_n) \in A) \le \limsup_{n \to \infty} \mathbb{P}(X_n \in \overline{\{h \in A\}}) \le \mathbb{P}(X \in \overline{\{h \in A\}}),$$

where the last inequality follows again from Portmanteau theorem (d) since $\overline{\{h \in A\}}$ is clearly closed and $X_n \stackrel{\text{w}}{\to} X$. Finally, as $\mathbb{P}(X \in C_h) = 1$,

$$\mathbb{P}(X \in \overline{\{h \in A\}}) = \mathbb{P}(X \in \overline{\{h \in A\}} \cap C_h) \leq \mathbb{P}(X \in \{h \in \overline{A}\})$$

from Proposition 2.2.1, i.e.,

$$\lim_{n\to\infty} \mathbb{P}(h(X_n)\in A) \le \mathbb{P}(X\in\{h\in\overline{A}\}) = \mathbb{P}(X\in\{h\in A\}) = \mathbb{P}(h(X)\in A)$$

since A is closed, hence we're done.

Example. Let d=1 and $X_n \stackrel{\text{w}}{\to} X$ where X is continuous. Then $1/X_n \stackrel{\text{w}}{\to} 1/X$ and $X_n^2 \stackrel{\text{w}}{\to} X^2$.

Proof. For the case of $X^2 \stackrel{\text{w}}{\to} X^2$, continuous mapping theorem clearly applies with $h(x) = x^2$. For the first case, consider

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

This means $C_h = \mathbb{R} \setminus \{0\}$. Then, we just need to show $\mathbb{P}(X \in C_h) = 1$ and apply continuous mapping theorem. Observe that this is the same as asking $\mathbb{P}(X = 0) = 0$, which is true when X is continuous.^a

 a Even if X is not continuous, as long as this is true we can conclude the same thing.

Another useful theorem for proving weak convergence is the following.

Theorem 2.2.3 (Converging together). Let $X_n \stackrel{\text{w}}{\to} X$, and if Y_n on the same probability space as X_n such that $||X_n - Y_n|| \stackrel{p}{\to} 0$, i.e., for all $\epsilon > 0$, $\mathbb{P}(||X_n - Y_n|| > \epsilon) \to 0$ as $n \to \infty$. Then, $Y_n \stackrel{\text{w}}{\to} X$.

We first see some applications.

Corollary 2.2.1. If $Y_n \stackrel{p}{\to} X$, then $Y_n \stackrel{\text{w}}{\to} X$. The converse holds as long as $\mathbb{P}(X = c) = 1$ for some constant c.

Proof. By considering $X_n = X$ for all n, Theorem 2.2.3 implies that if $Y_n \stackrel{p}{\to} X$, $Y_n \stackrel{\text{w}}{\to} X$. Conversely, if $Y_n \stackrel{\text{w}}{\to} c$, from Portmanteau theorem (c), for any fixed $\epsilon > 0$,

$$\underbrace{\mathbb{P}(c \in B(c, \epsilon))}_{1} \le \liminf_{n \to \infty} \mathbb{P}(Y_n \in B(c, \epsilon)),$$

implying $\mathbb{P}(Y_n \in B(c, \epsilon)) \to 1$, i.e., $\mathbb{P}(\|Y_n - c\| < \epsilon) \to 1$.

Lecture 5: Convergence in Distribution and Weak Convergence

Now we prove Theorem 2.2.3.

1 Feb. 9:30

Proof. From Portmanteau theorem (b), we want to prove that $\mathbb{E}\left[g(Y_n)\right] \to \mathbb{E}\left[g(X)\right]$ for all bounded and Lipschitz $g \colon \mathbb{R}^d \to \mathbb{R}$. Specifically, let $|g(x)| \leq C$ for all $x \in \mathbb{R}^d$ and $|g(x) - g(y)| \leq K||x - y||$ for all $x, y \in \mathbb{R}^d$. From triangle inequality,

$$|\mathbb{E}\left[g(Y_n)\right] - \mathbb{E}\left[g(X)\right]| \le |\mathbb{E}\left[g(Y_n)\right] - \mathbb{E}\left[g(X_n)\right]| + |\mathbb{E}\left[g(X_n)\right] - \mathbb{E}\left[g(X)\right]|.$$

Since $X_n \stackrel{\text{w}}{\to} X$, the second term goes to 0. As for the first term, since Y_n and X_n are in the same probability space, we see that

$$\begin{split} |\mathbb{E}\left[g(Y_n)\right] - \mathbb{E}\left[g(X_n)\right]| &= |\mathbb{E}\left[g(Y_n) - g(X_n)\right]| \\ &\leq \mathbb{E}\left[|g(Y_n) - g(X_n)|\right] \\ &= \mathbb{E}\left[|g(Y_n) - g(X_n)| \cdot \mathbb{1}_{\|X_n - Y_n\| > \epsilon}\right] + \mathbb{E}\left[|g(Y_n) - g(X_n)| \cdot \mathbb{1}_{\|X_n - Y_n\| \le \epsilon}\right] \\ &\leq 2C\mathbb{P}(\|X_n - Y_n\| > \epsilon) + K\epsilon\mathbb{P}(\|X_n - Y_n\| \le \epsilon) \\ &\leq 2C\mathbb{P}(\|X_n - Y_n\| > \epsilon) + K\epsilon. \end{split}$$

As $n \to \infty$, we finally have

$$\limsup_{n \to \infty} |\mathbb{E}\left[g(Y_n)\right] - \mathbb{E}\left[g(X)\right]| \le K\epsilon$$

for all $\epsilon > 0$, by letting $\epsilon \to 0$, we're done.

We can now apply Theorem 2.2.3 to prove something similar as we have seen before in the case of convergence in probability.

^aRecall that $B(c,\epsilon)$ is the open ball centered at c with radius ϵ .

As previously seen. $X_n \xrightarrow{p} X$ and $Y_n \xrightarrow{p} Y$ if and only if $(X_n, Y_n) \xrightarrow{p} (X, Y)$.

Now, in the case of weak convergence, from continuous mapping theorem, we see that if $(X_n, Y_n) \xrightarrow{w} (X, Y)$, then $X_n \xrightarrow{w} X$ and $Y_n \xrightarrow{w} Y$. However, the converse needs not be true.

Example. Consider a random variable X on $(\Omega, \mathscr{F}, \mathbb{P})$, and let $X_n = X$, $Y_n = -X$ for all $n \geq 1$. If $X \sim \mathcal{N}(0,1)$, we see that $\mathbb{P}(X \in A) = \mathbb{P}(-X \in A)$ for all $A \subseteq \mathbb{R}^d$, implying $X_n \overset{\text{w}}{\to} X$ and $Y_n \overset{\text{w}}{\to} X$. However, this does not imply $(X_n, Y_n) \overset{\text{w}}{\to} (X, X)$ since otherwise, by continuous mapping theorem, $X_n + Y_n \overset{\text{w}}{\to} X + X = 2X$, which is not true since $X_n + Y_n = 0$.

But in the case of Y is a constant, the converse is actually true, and the result is quite useful.

Theorem 2.2.4 (Slutsky's theorem). If $X_n \stackrel{\mathbb{W}}{\to} X$ in \mathbb{R}^d and $Y_n \stackrel{p}{\to} c$ in \mathbb{R}^m , $\stackrel{a}{\to}$ then $(X_n, Y_n) \stackrel{\mathbb{W}}{\to} (X, c)$.

^aRecall that from Corollary 2.2.1, for a constant c, weak convergence is equivalent to convergence in probability.

Proof. Firstly, we show that $(X_n, c) \xrightarrow{w} (X, c)$. Indeed, since for every continuous and bounded $g \colon \mathbb{R}^{d+m} \to \mathbb{R}, \mathbb{E}\left[g(X_n, c)\right] \to \mathbb{E}\left[g(X, c)\right]$ follows directly from $X_n \xrightarrow{w} X$ with $g(\cdot, c)$ being continuous and bounded.

Secondly, we show that $\|(X_n, Y_n) - (X_n, c)\| \stackrel{p}{\to} 0$. This is easy since

$$||(X_n, Y_n) - (X_n, c)|| \le ||X_n - X_n|| + ||Y_n - c|| = ||Y_n - c||,$$

which goes to 0 in probability as we wish. Combining both with Theorem 2.2.3 gives the result. ■

Revisiting the counter-example, we see that now it's not the case when Y is a constant.

Corollary 2.2.2. If $X_n \stackrel{\text{W}}{\to} X$ and $Y_n \stackrel{p}{\to} c$ in \mathbb{R}^d , $X_n \pm Y_n \stackrel{\text{W}}{\to} X \pm c$.

2.2.3 Convergence in Distribution

So far, the notions of convergence we have talked about applies to general probability space, which needs not to be in \mathbb{R}^d in general. However, traditionally, the case in \mathbb{R}^d is considered first.

Definition 2.2.4 (Converge in distribution). Let (X_n) and X be random variables in \mathbb{R}^d . Then (X_n) converges in distribution to X, denoted as $X_n \stackrel{D}{\to} X$, if for all $(t_1, \ldots, t_d) \in C_{F_X}$,

$$F_{X_n}(t_1,\ldots,t_d)\to F_X(t_1,\ldots,t_d).$$

Specifically, recall that

$$F_{X_n}(t_1,\ldots,t_d) = \mathbb{P}(X_n^i \le t_i, \forall 1 \le i \le d) = \mathbb{P}(X_n \in (-\infty,t_1] \times \cdots \times (-\infty,t_d]),$$

same for F_X . We now see why we mentioned \mathbb{R}^d specifically:

Intuition. There's a conical ordering available in \mathbb{R}^d to define F_X and F_{X_n} .

There are more remarks to be made here.

Remark. X_n and X (in \mathbb{R}^d) do not have to be on the same probability space.

So far, we have talked about

Remark. $X_n \stackrel{\mathrm{TV}}{\to} X$ implies $X_n \stackrel{D}{\to} X$.

Proof. Since $X_n \stackrel{\mathrm{TV}}{\to} X$ means $\mathbb{P}(X_n \in A) \to \mathbb{P}(X \in A)$ uniformly in A, but $X_n \stackrel{D}{\to} X$ only requires the above holds for A in the form of $(-\infty, t_1] \times \cdots \times (-\infty, t_d]$, which is weaker.

Remark (De Moivre). Let $X_n \sim \text{Bin}(n, p)$, then for every $t \in \mathbb{R}$, as $n \to \infty$,

$$\mathbb{P}\left(\frac{X_n - np}{\sqrt{np(1-p)}} \le t\right) \to \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-u^2/2} du =: \phi(u) = \mathbb{P}(Z \le t)$$

where $Z \sim \mathcal{N}(0, 1)$.

Remark (Polya's theorem). If F_X is continuous, then for all $t \in \mathbb{R}^d$,

$$X_n \stackrel{D}{\to} X \Leftrightarrow \sup_{t \in \mathbb{R}} |F_{X_n}(t) - F_X(t)| \to 0.$$

Problem. Why not defined for all $t \in \mathbb{R}^d$, rather than $t \in C_{F_X}$?

Answer. Consider for d=1 with $X=c\in\mathbb{R}$, i.e., F_X is the step function at c. To show $X_n\stackrel{D}{\to}c$, we don't have to show $\mathbb{P}(X_n\leq c)\to\mathbb{P}(X\leq c)=1$. Otherwise, if we need to show this for all t, in particular, c, $X_n=c+1/n$ would not satisfy this.

Remark. Let X_n and X be in \mathbb{Z} . Let f_n and f be their corresponding PMF's, then

$$f_n \to f \Leftrightarrow X_n \stackrel{\mathrm{TV}}{\to} X \Leftrightarrow X_n \stackrel{D}{\to} X.$$

Proof. The forward implications are clear, so we just need to show $X_n \stackrel{D}{\to} X$ implies $f_n \to f$. Since for every $t \in \mathbb{Z}$, since X_n and X are discrete in \mathbb{Z} ,

$$f_n(t) = \mathbb{P}(X_n = t) = \mathbb{P}(X_n \le t + \epsilon) - \mathbb{P}(X_n \le t - \epsilon)$$

for some $\epsilon > 0$ small enough. Now, as $t \pm \epsilon$ are in C_X clearly, $X_n \stackrel{D}{\to} X$ implies

$$\mathbb{P}(X_n \le t + \epsilon) \to \mathbb{P}(X \le t + \epsilon).$$

and the same holds for $t - \epsilon$, hence

$$f_n(t) = \mathbb{P}(X_n = t) = \mathbb{P}(X_n \le t + \epsilon) - \mathbb{P}(X_n \le t - \epsilon) \to \mathbb{P}(X \le t + \epsilon) - \mathbb{P}(X \le t - \epsilon) = \mathbb{P}(X = t) = f(t).$$

As this holds for every $t \in \mathbb{Z}$, we're done.

The virtue of convergence in distribution is that it's actually just a renamed version of weak convergence in \mathbb{R}^d .

Theorem 2.2.5. Given X_n and X in \mathbb{R}^d , then $X_n \stackrel{\text{w}}{\to} X$ if and only if $X_n \stackrel{D}{\to} X$.

Proof. We prove for the case of d=1, then it's easy to see the same holds for $d\geq 1$. For the forward direction, we want to show that for all $t\in C_{F_X}$,

$$\mathbb{P}(X_n < t) \to \mathbb{P}(X < t).$$

Note that $\mathbb{P}(X \leq t) = \mathbb{P}(X \in (-\infty, t])$ and $\mathbb{P}(X_n \leq t) = \mathbb{P}(X_n \in (-\infty, t])$, hence, from Portmanteau theorem (e) with $A = (-\infty, t]$, $X_n \stackrel{\text{w}}{\to} X$ is equivalently as saying $\mathbb{P}(X_n \leq t) \to \mathbb{P}(X \leq t)$ if

$$\mathbb{P}(X \in \partial(-\infty, t]) = \mathbb{P}(X \in \{t\}) = \mathbb{P}(X = t)$$

is 0. This is indeed the case since $t \in C_{F_X}$, hence we're done.

To show the backward direction, we need the following lemma.

Lemma 2.2.3. $X_n \stackrel{D}{\to} X$ if and only if for all $x \in \mathbb{R}^d$,

$$F_X(x^-) \le \liminf_{n \to \infty} F_{X_n}(x^-) \le \liminf_{n \to \infty} F_{X_n}(x) \le \limsup_{n \to \infty} F_{X_n}(x) \le F_X(x).$$

Proof. The backward direction is clear, so we prove the forward direction. When $x \in C_{F_X}$, we're clearly done, so consider $x \notin C_{F_X}$. Firstly, note that $|C_{F_X}^c|$ is countable, so there exists $(x_k) \nearrow x$ and $(y_k) \searrow x$, both in C_{F_X} . Hence, for all $n \ge 1$ and $k \ge 1$,

$$F_{X_n}(x_k) \le F_{X_n}(x) \le F_{X_n}(y_k)$$

as F_{X_n} is increasing. We now have for every $k \geq 1$,

$$\begin{split} F_X(x_k) &= \lim_{n \to \infty} F_{X_n}(x_k) & x_k \in C_{F_X} \\ &\leq \liminf_{n \to \infty} F_{X_n}(x^-) \\ &\leq \liminf_{n \to \infty} F_{X_n}(x) & F_{X_n} \text{ is increasing} \\ &\leq \limsup_{n \to \infty} F_{X_n}(x) \\ &\leq \limsup_{n \to \infty} F_{X_n}(y_k) = F_X(y_k). & y_k \in C_{F_X} \end{split}$$

By taking $k \to \infty$, $F_X(x_k) \to F_X(x^-)$, while $F_X(y_k) \to F_X(x)$, and we're done.

The proof will be continued...

Lecture 6: Stochastic Boundedness and Delta Theorem

Before we finish the proof of Theorem 2.2.5, we need recall one important characterization of liminf.

2 Feb. 17:30

As previously seen. Given two real sequence x_n and y_n ,

$$\liminf_{n \to \infty} (x_n + y_n) \ge \liminf_{n \to \infty} x_n + \liminf_{n \to \infty} y_n,$$

where the equality holds when either x_n or y_n converges (not if and only if).

We can then finish the proof of Theorem 2.2.5.

Proof of Theorem 2.2.5 (cont.) Now we can prove the backward direction. Form Portmanteau theorem (c), it suffices to show that for every open $A \subseteq \mathbb{R}$, we have

$$\mathbb{P}(X \in A) \le \liminf_{n \to \infty} \mathbb{P}(X_n \in A).$$

From the elementary analysis, we see that it suffices to show when A = (a, b) since when $A \subseteq \mathbb{R}$ is open, one can write $A = \bigcup_{k=1}^{\infty} (a_k, b_k)$ where (a_k, b_k) 's disjoint, and have

$$\mathbb{P}(X \in A) = \sum_{k=1}^{\infty} \mathbb{P}(X \in (a_k, b_k))$$

$$\leq \sum_{k=1}^{\infty} \liminf_{n \to \infty} \mathbb{P}(X_n \in (a_k, b_k))$$
 assume true for intervals
$$\leq \liminf_{n \to \infty} \sum_{k=1}^{\infty} \mathbb{P}(X_n \in (a_k, b_k)) = \liminf_{n \to \infty} \mathbb{P}(X_n \in A),$$

where the last inequality follows from induction on $\liminf_{n\to\infty} (x_n+y_n) \ge \liminf_{n\to\infty} x_n + \liminf_{n\to\infty} y_n$. Now, we show that $\mathbb{P}(X \in A) \le \liminf_{n\to\infty} \mathbb{P}(X_n \in A)$ when A = (a,b).

 $[^]a$ Recall that the distribution function is always right-continuous.

Claim. $\mathbb{P}(X \in (a,b)) \leq \liminf_{n \to \infty} \mathbb{P}(X_n \in (a,b)).$

Proof. Observe that $\mathbb{P}(X \in (a,b)) = F_X(b^-) - F_X(a)$, with Lemma 2.2.3, we further have

$$\begin{split} \mathbb{P}(X \in (a,b)) &= F_X(b^-) - F_X(a) \\ &\leq \liminf_{n \to \infty} F_{X_n}(b^-) - \left(\limsup_{n \to \infty} F_{X_n}(a)\right) \\ &\leq \liminf_{n \to \infty} F_{X_n}(b^-) + \liminf_{n \to \infty} (-F_{X_n}(a)) \\ &\leq \liminf_{n \to \infty} \left(F_{X_n}(b^-) - F_{X_n}(a)\right) = \liminf_{n \to \infty} \mathbb{P}(X_n \in (a,b)), \end{split}$$

which proves the claim.

This proves the case of d = 1.

Theorem 2.2.5 means that when talking about random vectors, we can use every result we have proved for the case of weak convergence. Let's see one application, which uses weak convergence's result but now prove something about the distribution.

Proposition 2.2.2. If $X_n \stackrel{D}{\to} X$ and $t_n \to t \in C_{F_X}$, then $\mathbb{P}(X_n \le t_n) \to \mathbb{P}(X \le t)$.

Proof. We see that from Corollary 2.2.2, $X_n - t_n \stackrel{\text{w}}{\to} X - t$, i.e., $X_n - t_n \stackrel{D}{\to} X - t$. Hence,

$$\mathbb{P}(X_n \le t_n) = \mathbb{P}(X_n - t_n \le 0) = F_{X_n - t_n}(0) \to F_{X - t}(0) = \mathbb{P}(X - t \le 0)$$

as long as $0 \in C_{F_{X-t}}$, i.e., $\mathbb{P}(X-t=0) = \mathbb{P}(X=t) = 0$, which is just $t \in C_{F_X}$ as we assumed.

2.3 Stochatsic Boundedness

So far we have been talking about the notion of convergence, now we switch the gear a bit and consider boundedness. In this section, let $(X_i)_{i\in I}$ be a family of d-dimensional random vectors with the index set I, which can be either finite or infinite.

Definition 2.3.1 (Bounded in probability). $(X_i)_{i \in I}$ is said to be bounded in probability if for every $\epsilon > 0$, there exists an M > 0 such that for every $i \in I$,

$$\mathbb{P}(\|X_i\| \ge M) < \epsilon.$$

In other words, for every $\epsilon > 0$, there exists an M > 0 such that $\mathbb{P}(\|X_i\| < M) \ge 1 - \epsilon$ for every $i \in I$.

Intuition. For any arbitrary large probability close to 1 we want, one can find an upper-bound M on $||X_i||$ uniformly for all $i \in I$.

Note. When $X_i = X$ for every $i \in I$, $(X_i)_{i \in I}$ is trivially bounded in probability.

Proof. Since if not, there exists $\epsilon > 0$, for every M > 0, $\mathbb{P}(\|X\| \ge M) \ge \epsilon$. Then as $M \to \infty$, $\mathbb{P}(\|X\| = \infty) \ge \epsilon$, which is a contradiction since $\|X\| = \infty$.

Remark. When I is finite, $(X_i)_{i \in I}$ is also trivially bounded in probability. On the other hand, when I is infinite, by considering $X_n = n$ (deterministic), which is not bounded in probability anymore.

Remark. If $(X_i)_{i\in I}$ is bounded in L^p for some p>0, i.e., $\sup_{i\in I} \mathbb{E}[\|X_i\|^p] < \infty$, then $(X_i)_{i\in I}$ is bounded in probability.

Proof. Since for any $\epsilon > 0$, from Markov's inequality,

$$\mathbb{P}(\|X_i\| > M) \le \frac{\mathbb{E}\left[\|X_i\|^p\right]}{M^p},$$

which can be made less than ϵ since $\sup_{i \in I} \mathbb{E}[\|X_i\|^p] < \infty$, for M large enough it'll be satisfied. \circledast

2.3.1 Convergence and Boundedness

Recall the following fact in elementary analysis.

As previously seen. If a deterministic sequence in \mathbb{R} converges, then it's bounded.

In our context, we might expect something like "if $X_n \stackrel{p}{\to} X$, then (X_n) is bounded in probability." In fact, we have the following "stronger" result where we only require convergence in distribution.

Proposition 2.3.1. If $X_n \stackrel{D}{\to} X$, then (X_n) is bounded in probability.

Proof. Fix an $\epsilon > 0$. There is an M > 0 such that $\mathbb{P}(\|X\| \ge M) < \epsilon$ since this is a single random vector. To relate this back to X_n , from Portmanteau theorem (d),

$$\epsilon > \mathbb{P}(\|X\| \ge M) = \mathbb{P}(X \in B^c(0, M)) \ge \limsup_{n \to \infty} \mathbb{P}(X_n \in B^c(0, M)) = \limsup_{n \to \infty} \mathbb{P}(\|X_n\| > M).$$

In other words,

$$\liminf_{n\to\infty} \mathbb{P}(\|X_n\| \le M) > 1 - \epsilon,$$

i.e., there exists an n_0 such that for every $n \ge n_0$, $\mathbb{P}(\|X_n\| \le M) \ge 1 - \epsilon$. As for those $n < n_0$, we can also find M' > 0 such that $\mathbb{P}(\|X_n\| \le M') > 1 - \epsilon$ for every $n < n_0$. Finally, by considering $M'' := \max(M, M')$, we have $\mathbb{P}(\|X_n\| \le M'') > 1 - \epsilon$, i.e., $\mathbb{P}(\|X_n\| > M) < \epsilon$ as desired.

There is a kind of converse theorem holds called Prokhorov's theorem, but we won't prove it here. Another useful characterization that generalizes our intuition in \mathbb{R} is the following.

Proposition 2.3.2. When d=1, if $X_n \stackrel{p}{\to} 0$ and Y_n is bounded in probability, then $X_n Y_n \stackrel{p}{\to} 0$.

Proof. Fix an $\epsilon > 0$. We want to show that $\mathbb{P}(|X_n Y_n| > \epsilon) \to 0$. This is because

$$\begin{split} \mathbb{P}(|X_nY_n| > \epsilon) &= \mathbb{P}(|X_nY_n| > \epsilon, |Y_n| > M) + \mathbb{P}(|X_nY_n| > \epsilon, |Y_n| \leq M) \\ &\leq \mathbb{P}(|Y_n| > M) + \mathbb{P}(|X_nY_n| > \epsilon, |Y_n| \leq M) \leq \mathbb{P}(|Y_n| > M) + \mathbb{P}(|X_n| > \epsilon/M) \end{split}$$

for any M. Now, we see that

- since Y_n is bounded in probability, there's an M > 0 such that $\mathbb{P}(|Y_n| > M) < \epsilon$ for all n;
- since $X_n \xrightarrow{p} 0$, for the M (depends on the fixed ϵ) above, $\mathbb{P}(|X_n| > \epsilon/M) \to 0$ as $n \to \infty$.

We see that the second term always goes to 0, while the first term can always be upper-bounded by ϵ . Hence, by letting $\epsilon \to 0$, we're done.

The analogy to the case in \mathbb{R} is the following.

Intuition. In \mathbb{R} , if $a_n \to 0$ and b_n is bounded, $a_n b_n \to 0$.

We often write the following.

Notation. We write $X_n = o_p(1)$ for $X_n \stackrel{p}{\to} 0$, and $X_n = O_p(1)$ when (X_n) is bounded in probability.

Let's see one important application. Consider an estimator T_n of θ , and a deterministic sequence b_n which goes to ∞ . In this case, we often have

$$b_n(T_n-\theta) \stackrel{D}{\to} Y.$$

Example. When $X_n \sim \text{Bin}(n, p)$, then

$$\frac{\sum_{i=1}^{n} X_i - np}{\sqrt{np(1-p)}} = \sqrt{\frac{n}{p(1-p)}} \left(\frac{\sum_{i=1}^{n} X_i}{n} - p \right) \to Y \sim \mathcal{N}(0,1)$$

This allows us to compute the rate of convergence and the limiting distribution. But what can we say when we care about $g(T_n)$ for a function g?

Theorem 2.3.1 (Delta method). Let (T_n) be d-dimensional random vectors with a deterministic sequence $b_n \to \infty$ such that $b_n(T_n - \theta) \stackrel{D}{\to} Y$. If $g: \mathbb{R}^d \to \mathbb{R}^m$ is differentiable at θ , then

$$||b_n(g(T_n) - g(\theta)) - \nabla g(\theta)b_n(T_n - \theta)|| \stackrel{p}{\to} 0.$$

Proof. Since g is differentiable at θ , as $x \to \theta$,

$$\frac{g(x) - g(\theta) - \nabla g(\theta)(x - \theta)}{\|x - \theta\|} \to 0.$$

Let $r(x) := g(x) - g(\theta) - \nabla g(\theta)(x - \theta)$ to be the remainder, and consider

$$h(x) = \begin{cases} 0, & \text{if } x = \theta; \\ \frac{r(x)}{\|x - \theta\|}, & \text{if } x \neq \theta, \end{cases}$$

which is continuous at θ . Rewriting everything, we have

$$r(x) = g(x) - g(\theta) - \nabla g(\theta)(x - \theta) = h(x)||x - \theta||$$

for every $x \in \mathbb{R}^d$. Now, let $x = T_n$, multiply both sides by b_n , and take the norm, we see that

$$||b_n(g(T_n) - g(\theta)) - \nabla g(\theta)b_n(T_n - \theta)|| = ||h(T_n)|| ||b_n(T_n - \theta)||.$$

We want to show that the right-hand sides goes to 0 in probability. Observe that it's enough to show $||h(T_n)|| = o_p(1)$ and $||b_n(T_n - \theta)|| \in O_p(1)$. Indeed:

- $||b_n(T_n \theta)|| \in O_p(1)$: Since $b_n(T_n \theta) \stackrel{D}{\to} Y$, with continuous mapping theorem and the fact that $||\cdot||$ is continuous, $||b_n(T_n \theta)|| \stackrel{D}{\to} ||Y||$, so $||b_n(T_n \theta)|| \in O_p(1)$ by Proposition 2.3.1.
- $||h(T_n)|| = o_p(1)$: since $b_n \to \infty$,

$$||T_n - \theta|| = \frac{1}{b_n} ||b_n(T_n - \theta)|| \stackrel{p}{\to} 0,$$

i.e., $T_n \stackrel{p}{\to} \theta$. This implies $h(T_n) \stackrel{p}{\to} h(\theta)$ again by continuous mapping theorem with h being continuous at θ . This further implies $||h(T_n)|| \stackrel{p}{\to} 0$ as we desired.

Combining the above, we prove the result.

Hence, we see that the answer to our original question is rather simple: as $b_n(T_n - \theta) \stackrel{D}{\to} Y$,

$$b_n(g(T_n) - g(\theta)) \stackrel{D}{\to} \nabla g(\theta) \cdot Y$$

for any differentiable q at θ .

^aThis involves continuous mapping theorem and Corollary 2.2.1 since $h(\theta) = 0$, a constant (so does its norm).

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: geloxygtxlec.
- [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).